
Sparse Kernel Machines

*Prof. Nicholas Zabaras
Center for Informatics and Computational Science*

<https://cics.nd.edu/>

*University of Notre Dame
Notre Dame, Indiana, USA*

Email: nzabaras@gmail.com

URL: <https://www.zabaras.com/>

March 12, 2019

SVM - Contents

- ❑ [SVM for Regression](#), [SVM for Classification](#), [Linear Classifiers](#), [Classifier Margin](#), [Support Vectors](#), [Computing the margin width](#)
- ❑ [SVM for the non-separable case](#)
- ❑ [Probabilistic Interpretation of SVM](#), [SVM Vs. Logistic Regression](#)
- ❑ [SVM for the Separable Case](#), [Kernel Trick and Non-Linear Classifier](#), [Regularization and Kernel Parameter Estimation](#), [Probabilistic Response](#), [Multiclass classification](#)

Following closely:

Bishop CM, [Pattern Recognition and Machine Learning](#), Springer, 2006 (Chapter 7)

Murphy, K. [Machine Learning: A probabilistic Approach](#) (Chapter 14)

Machine Learning, University of Notre Dame, Notre Dame, IN, USA (Spring 2019, N. Zabarás)

RVM - Contents

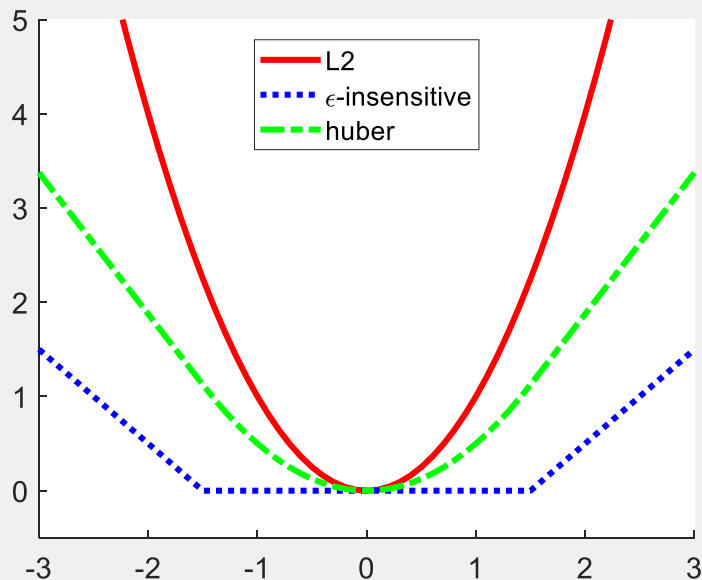
- ❑ Relevance Vector Machines: Likelihood, Prior, Posterior, Predictive Distribution
- ❑ Evidence Approximation, Evidence, Two useful matrix identities
- ❑ Fast training algorithm
- ❑ Splitting the evidence, Analysis of the evidence
- ❑ Complete algorithm, Example
- ❑ Fast updates
- ❑ Extension to Multiple Outputs, Numerical Stability Concerns

Support Vector Machines

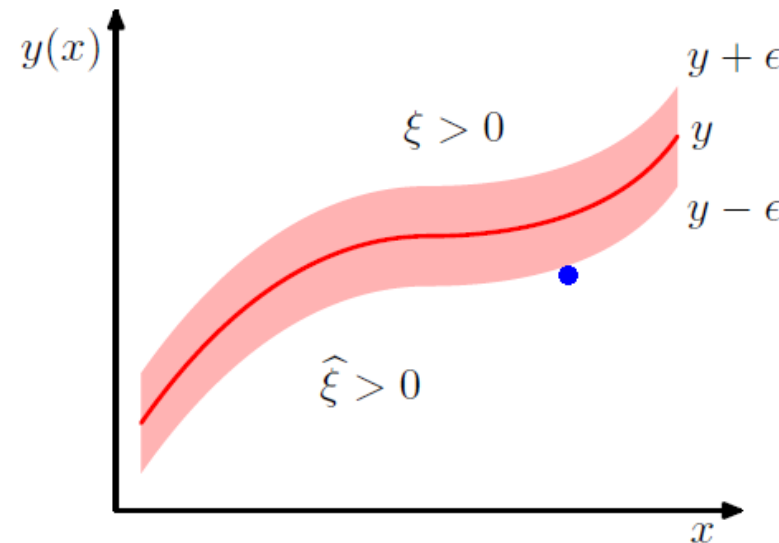
SVMs for Regression

- ❑ In our approach to regression, \mathbf{w} depends on all training inputs. We are here looking for sparse approximations.
- ❑ Consider the ϵ –insensitive loss function: $L_\epsilon(y, \hat{y}) \equiv \begin{cases} 0 & \text{if } |y - \hat{y}| < \epsilon \\ |y - \hat{y}| - \epsilon & \text{if } |y - \hat{y}| > \epsilon \end{cases}$
leading to the following objective function:

$$J = C \sum_{i=1}^N L_\epsilon(y_i, \hat{y}_i) + \frac{1}{2} \|\mathbf{w}\|^2, \hat{y}_i = f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i + w_0, C = 1/\lambda$$



[huberLossDemo.m](#)
from [PMTK3](#)



SVMs for Regression

- In our approach to regression, w depends on all training inputs. We are here looking for sparse approximations.

- Consider the ϵ –insensitive loss function: $L_\epsilon(y, \hat{y}) \equiv \begin{cases} 0 & \text{if } |y - \hat{y}| < \epsilon \\ |y - \hat{y}| - \epsilon & \text{if } |y - \hat{y}| > \epsilon \end{cases}$
leading to the following objective function:

$$J = C \sum_{i=1}^N L_\epsilon(y_i, \hat{y}_i) + \frac{1}{2} \|\mathbf{w}\|^2, \hat{y}_i = f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i + w_0, C = 1/\lambda$$

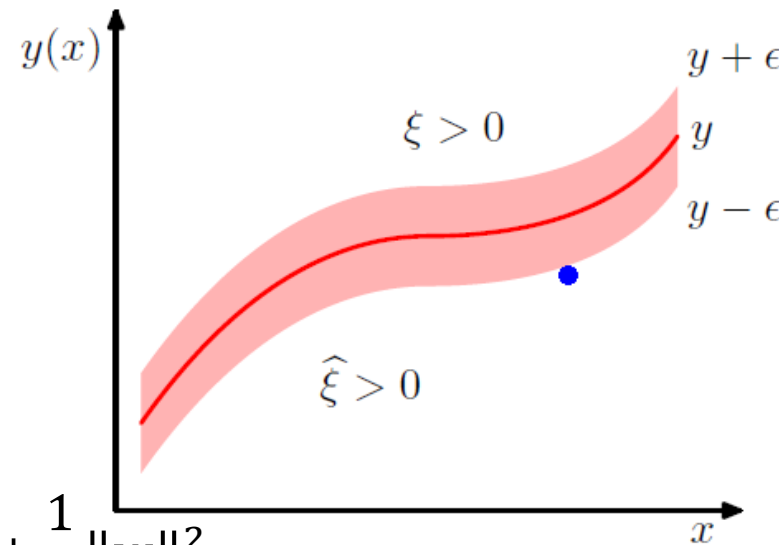
- Introduce slack variables to describe the degree to which each point lies outside the ϵ – tube:

$$y_i \leq f(\mathbf{x}_i) + \epsilon + \xi_i^+, \xi_i^+ \geq 0, i = 1, \dots, n$$

$$y_i \geq f(\mathbf{x}_i) - \epsilon - \xi_i^-, \xi_i^- \geq 0, i = 1, \dots, n$$

$$\min_{w, \xi_i^+ \geq 0, \xi_i^- \geq 0} J = \min_{w, \xi_i^+ \geq 0, \xi_i^- \geq 0} C \sum_{i=1}^N (\xi_i^+ + \xi_i^-) + \frac{1}{2} \|\mathbf{w}\|^2$$

Quadratic optimization
 $2N + D + 1$ variables



SVMs for Regression

- ❑ One can show that: $\hat{\mathbf{w}} = \sum_{i=1}^N \alpha_i \mathbf{x}_i$, $\alpha_i \geq 0$.
- ❑ The α vector is sparse since we do not care about errors $\leq \varepsilon$. The \mathbf{x}_i 's are called support vectors – these are points for which the errors lie on or outside the ε tube.
- ❑ Making predictions:

$$\hat{y}(\mathbf{x}) = \hat{w}_0 + \hat{\mathbf{w}}^T \mathbf{x} = \hat{w}_0 + \sum_{i=1}^N \alpha_i \mathbf{x}_i^T \mathbf{x} = \hat{w}_0 + \sum_{i=1}^N \alpha_i k(\mathbf{x}_i, \mathbf{x})$$

- ❑ In the last step above we used the kernel trick to substitute $\mathbf{x}_i^T \mathbf{x}$ with $k(\mathbf{x}_i, \mathbf{x})$.
- ❑ The kernel trick is needed to prevent underfitting (ensure the feature vector is sufficient rich that a linear classifier can separate the data).
- ❑ The sparsity and large margin are needed to prevent overfitting.
- Schoelkopf, B. and A. Smola (2002). [Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond](#). MIT Press.

SVM for Classification

- We have seen that the negative log-likelihood for the logistic regression model is

$$L_{nll}(y, \eta) = -\log p(y|\mathbf{x}, \mathbf{w}) = \log(1 + e^{-y\eta})$$

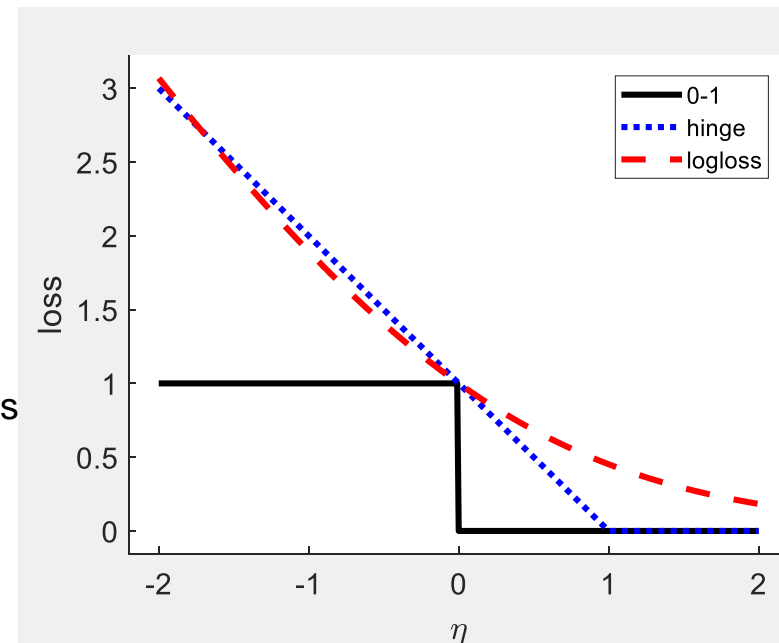
- It is a convex upper bound on the 0 – 1 risk of a binary classifier, where $\eta = f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ is the log odds ratio, and we have assumed the labels are $y \in \{1, -1\}$ rather than $\{0, 1\}$.

- We can replace the NLL loss with the **hinge loss**

$$L_{hinge}(y, \eta) = \max(0, 1 - y\eta) = (1 - y\eta)_+$$

The horizontal axis is the margin $y\eta$, the vertical axis is the loss. The log loss uses log base 2.

[hingeLossPlot.m](#)
from [PMTK3](#)



SVM for Classification

- Our optimization problem is as follows:

$$\min_{\mathbf{w}, w_0} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N (1 - y_i f(\mathbf{x}_i))_+$$

- Introducing slack variables we can write this as a quadratic program in $N + D + 1$ variables:

$$\min_{\mathbf{w}, w_0, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i, \quad \xi_i \geq 0, y_i(\mathbf{x}_i^T \mathbf{w} + w_0) \geq 1 - \xi_i, i = 1, \dots, N$$

- We can eliminate \mathbf{w} , w_0 and ξ_i , and just solve the N dual variables, which correspond to the Lagrange multipliers for the constraints.

Standard solvers take $\mathcal{O}(N^3)$ time. **Sequential minimal optimization** takes $\mathcal{O}(N^2)$. Linear SVMs, take $\mathcal{O}(N)$ time to train.

- Platt, J. (1998). [Using analytic QP and sparseness to speed training of support vector machines](#). In *NIPS*.
- Joachims, T. (2006). [Training Linear SVMs in Linear Time](#). In *Proc. Of the Int'l Conf. on Knowledge Discovery and Data Mining*.
- Bottou, L., O. Chapelle, D. DeCoste, and J. Weston (Eds.) (2007). [Large Scale Kernel Machines](#). MIT Press.

SVMs for Classification

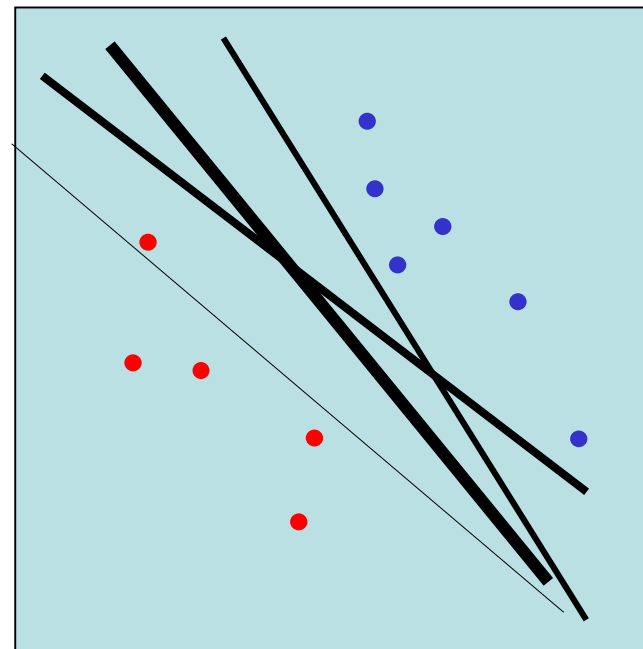
- ❑ One can show that: $\hat{\mathbf{w}} = \sum_{i=1}^N \alpha_i \mathbf{x}_i$, $\alpha_i = \lambda_i y_i$, α is sparse
- ❑ The \mathbf{x}_i 's for which $\alpha_i > 0$ are called support vectors – these are points which are incorrectly classified or are classified correctly but are on or inside the margin.
- ❑ Making predictions:

$$\hat{y}(\mathbf{x}) = \text{sgn}(\hat{\mathbf{w}}_0 + \hat{\mathbf{w}}^T \mathbf{x}) = \text{sgn}(\hat{\mathbf{w}}_0 + \sum_{i=1}^N \alpha_i k(\mathbf{x}_i, \mathbf{x}))$$

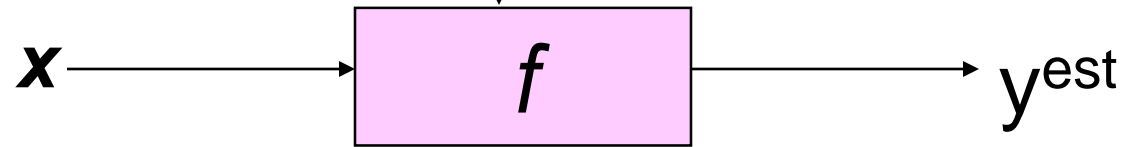
- ❑ It takes $\mathcal{O}(sD)$ time to compute where $s \leq N$ is the number of support vectors.

Linear Classifiers

- Suppose we use a big set of features to ensure that the two classes are linearly separable. What is the best separating line to use?
- ✓ The Bayesian answer is to use them all (including ones that do not quite separate the data.)
- ✓ Weight each line by its posterior probability (i.e. by a combination of how well it fits the data and how well it fits the prior).
- ✓ Is there an efficient way to approximate the correct Bayesian answer?

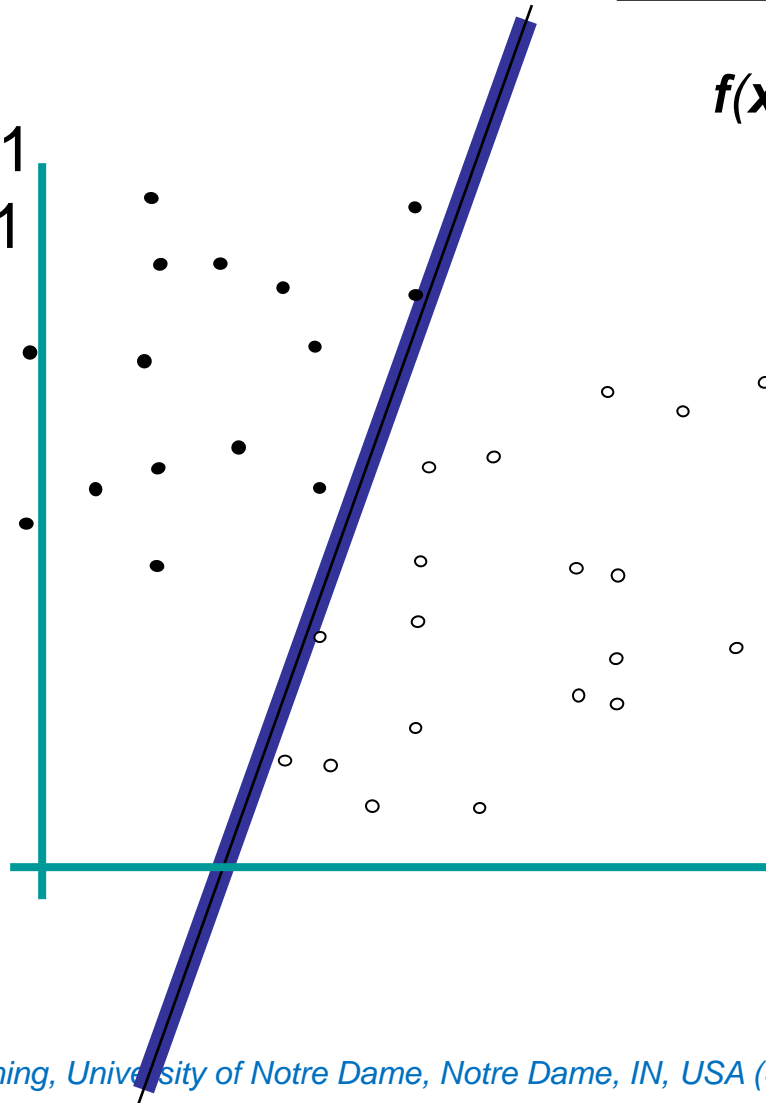


Classifier Margin



$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

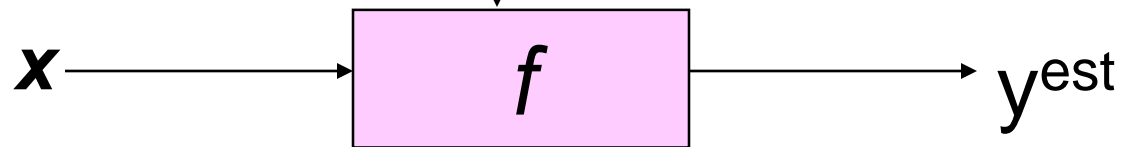
- denotes +1
- denotes -1



Linear SVM

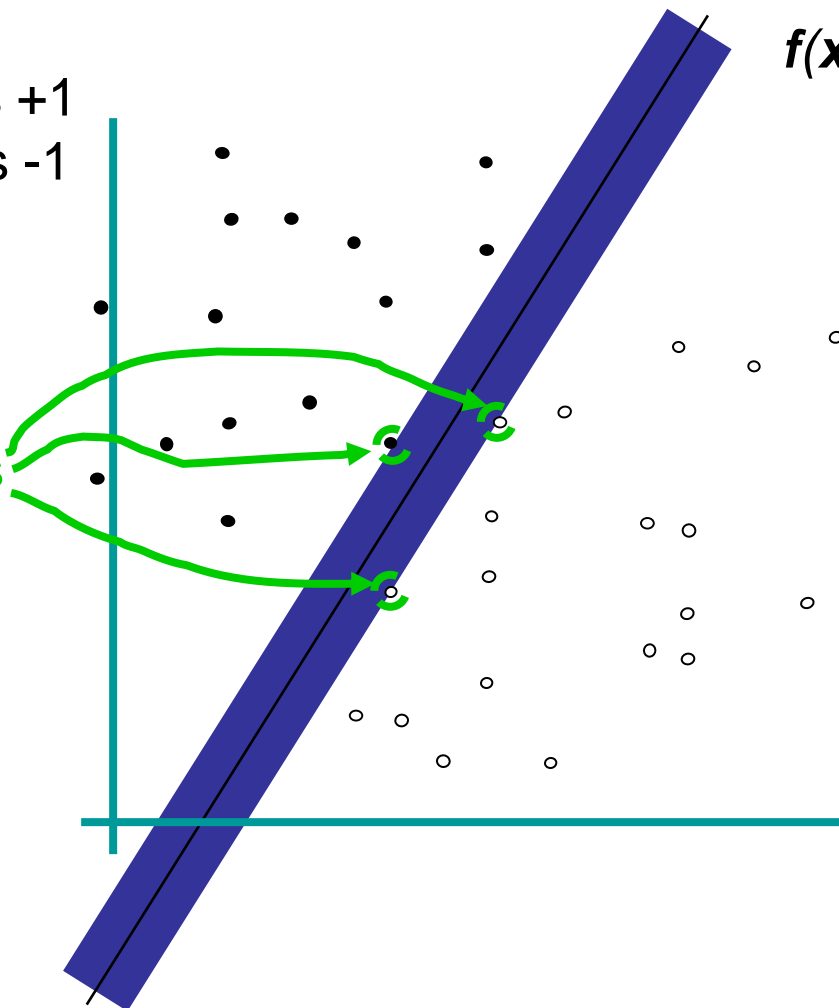
Define the **margin** of a linear classifier as the width that the boundary could be increased by before hitting a datapoint.

Support Vectors



- denotes +1
- denotes -1

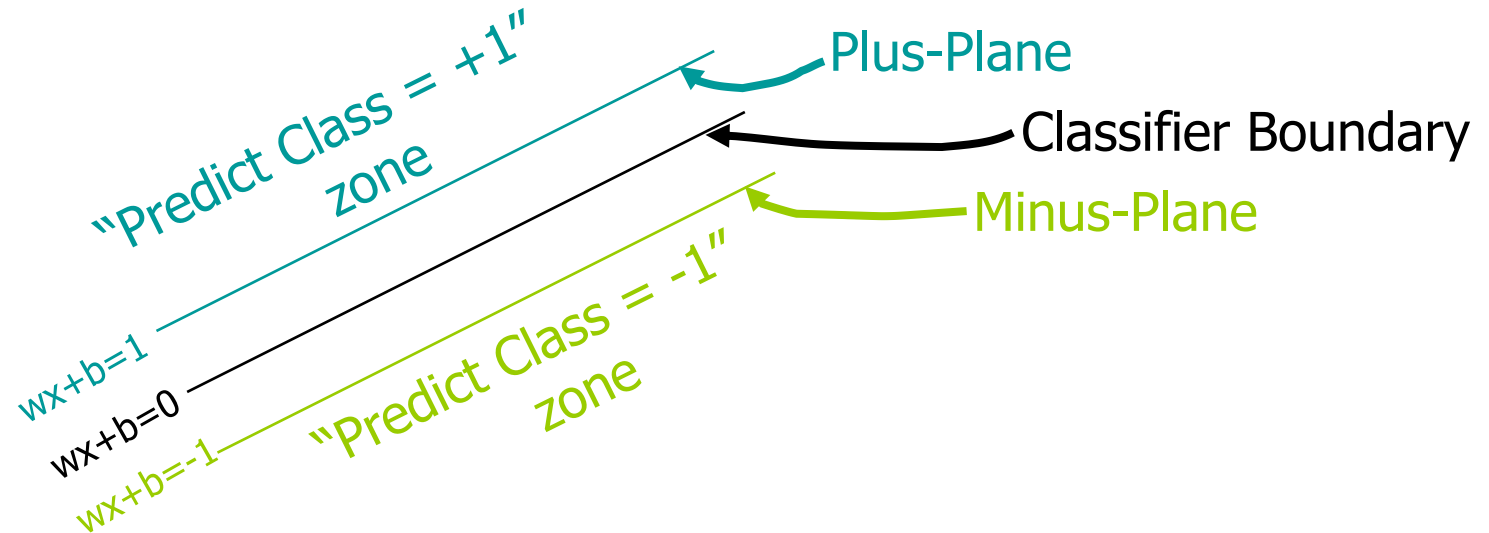
Support Vectors
are those
datapoints that
the margin
pushes up
against



$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

The maximum
margin linear
classifier is the
linear classifier
with the maximum
margin.

Specifying a line and margin



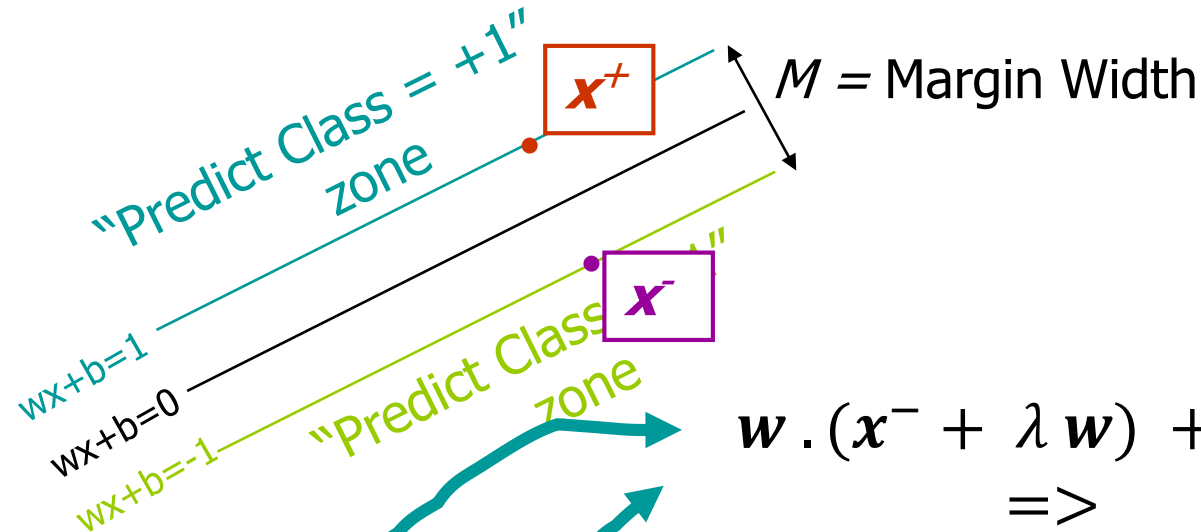
□ Plus-plane = $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = +1 \}$

□ Minus-plane = $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = -1 \}$

Classify as.. +1 if $\mathbf{w} \cdot \mathbf{x} + b \geq 1$

-1 if $\mathbf{w} \cdot \mathbf{x} + b \leq -1$

Computing the margin width



What we know:

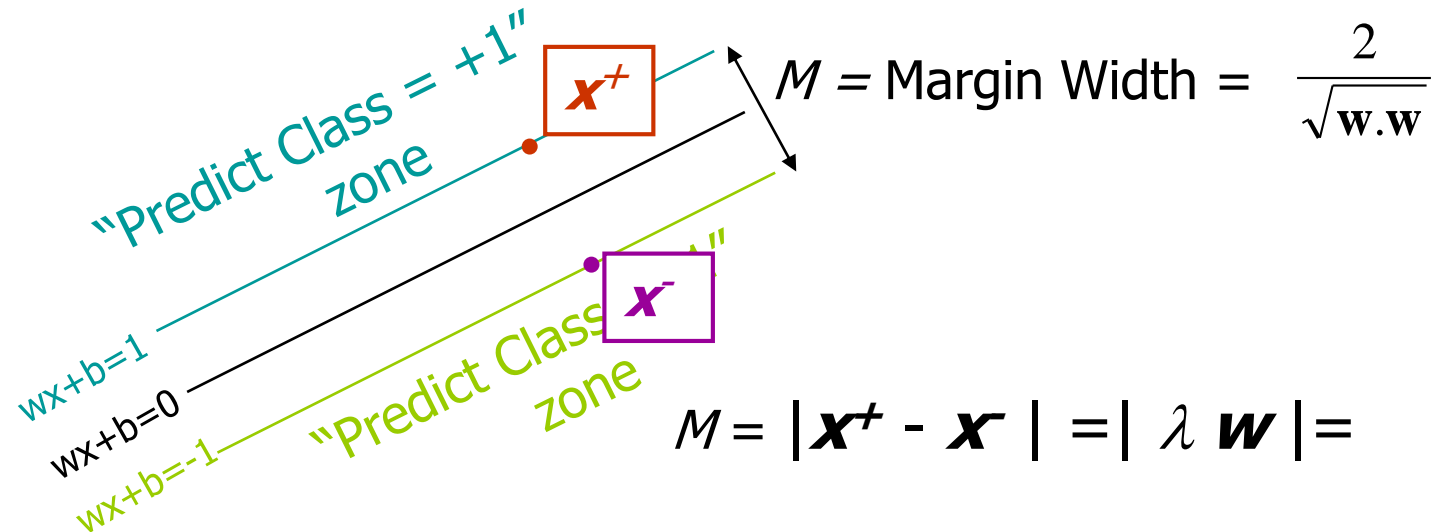
- $w \cdot x^+ + b = +1$
- $w \cdot x^- + b = -1$
- $x^+ = x^- + \lambda w$
- $|x^+ - x^-| = M$

It's now easy to get M
in terms of w and b

$$\begin{aligned} w \cdot (x^- + \lambda w) + b &= 1 \\ \Rightarrow \\ w \cdot x^- + b + \lambda w \cdot w &= 1 \\ \Rightarrow \\ -1 + \lambda w \cdot w &= 1 \\ \Rightarrow \end{aligned}$$

$$\lambda = \frac{2}{w \cdot w}$$

Computing the margin width



$$M = |\mathbf{x}^+ - \mathbf{x}^-| = |\lambda \mathbf{w}| =$$

$$= \lambda |\mathbf{w}| = \lambda \sqrt{\mathbf{w} \cdot \mathbf{w}}$$

$$= \frac{2\sqrt{\mathbf{w} \cdot \mathbf{w}}}{\mathbf{w} \cdot \mathbf{w}} = \frac{2}{\sqrt{\mathbf{w} \cdot \mathbf{w}}}$$

What we know:

$$\square \mathbf{w} \cdot \mathbf{x}^+ + b = +1$$

$$\square \mathbf{w} \cdot \mathbf{x}^- + b = -1$$

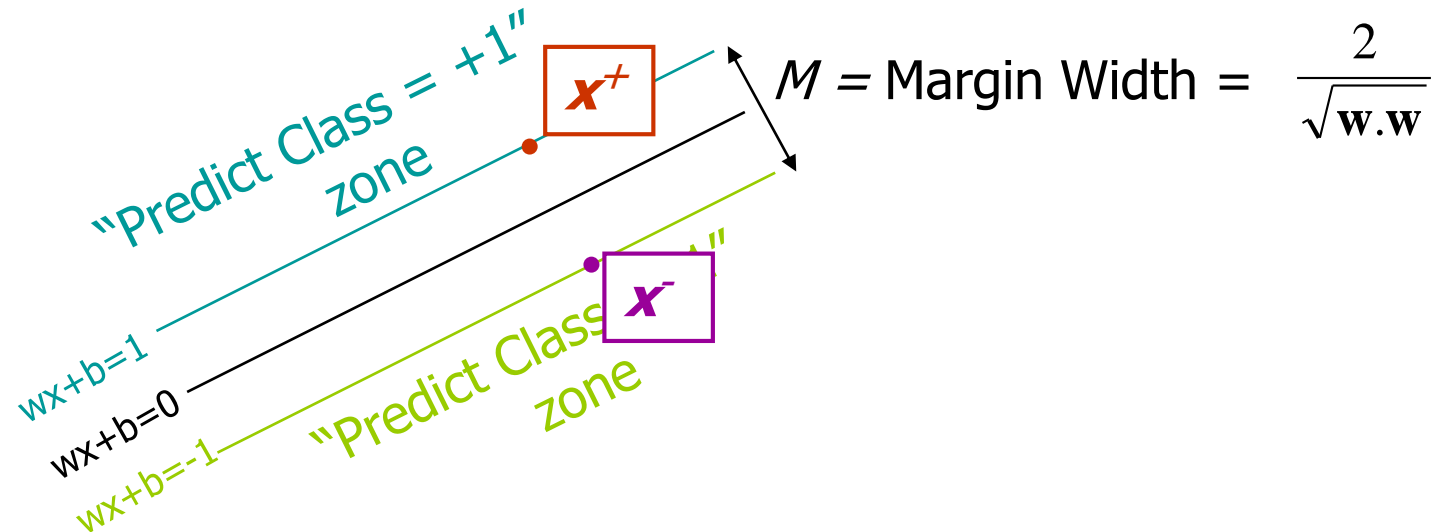
$$\square \mathbf{x}^+ = \mathbf{x}^- + \lambda \mathbf{w}$$

$$\square |\mathbf{x}^+ - \mathbf{x}^-| = M$$

\square

$$\lambda = \frac{2}{\mathbf{w} \cdot \mathbf{w}}$$

Learning the Maximum Margin Classifier



Given a guess of \mathbf{w} and b we can

- ☐ Compute whether all data points in the correct half-planes
- ☐ Compute the width of the margin

So now we just need to write a program to search the space of \mathbf{w} 's and b 's to find the widest margin that matches all the datapoints. *How?*

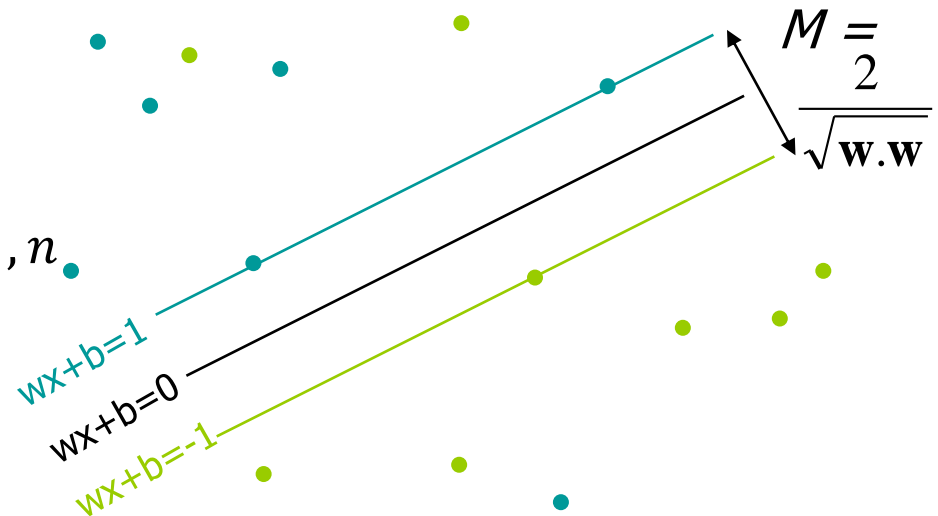
Support Vector Machine

- When the training examples are linearly separable we can maximize a geometric notion of margin (distance to the boundary) by minimizing the regularization penalty

$$\frac{1}{2} \|\mathbf{w}\|^2 = \frac{1}{2} \sum_{i=1}^D w_i^2$$

subject to the classification constraints

$$y_i[w_0 + \mathbf{x}_i^T \mathbf{w}] - 1 \geq 0, i = 1, \dots, n$$



- The solution is defined only on the basis of a subset of examples or “support vectors”

SVM for the non-separable case

- When the examples are not linearly separable we can modify the optimization problem slightly to add a penalty for violating the classification constraints: We minimize

$$\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i$$

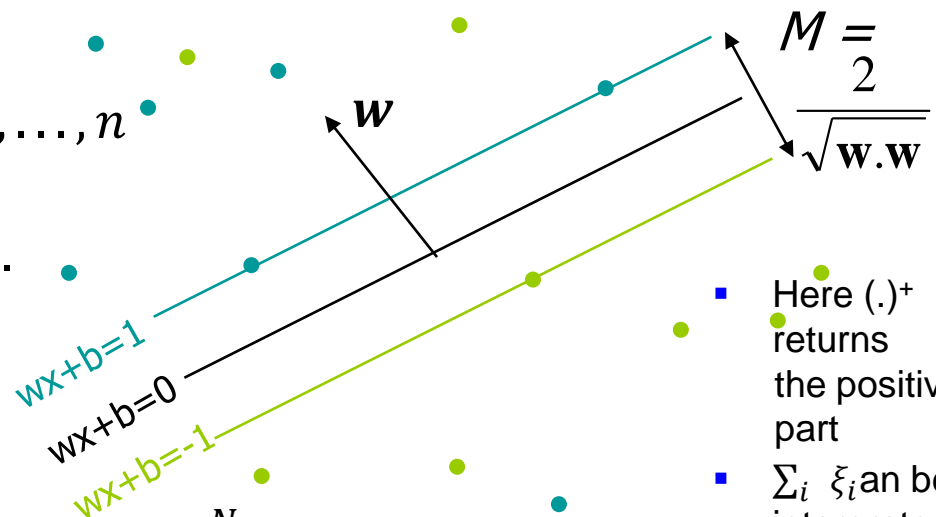
subject to the classification constraints

$$y_i [w_0 + \mathbf{x}_i^T \mathbf{w}] - 1 + \xi_i \geq 0, i = 1, \dots, n$$

- Here $\xi_i \geq 0$ are slack variables.

- Can write the minimization objective as:

$$C \sum_{i=1}^N (1 - y_i [w_0 + \mathbf{x}_i^T \mathbf{w}])^+ + \frac{1}{2} \|\mathbf{w}\|^2 = C \sum_{i=1}^N \xi_i + \frac{1}{2} \|\mathbf{w}\|^2$$



- Here $(.)^+$ returns the positive part
- $\sum_i \xi_i$ can be interpreted as the number of misclassified points

SVM for the non-separable case

- Note that the minimization objective

$$C \sum_{i=1}^N (1 - y_i [w_0 + \mathbf{x}_i^T \mathbf{w}])^+ + \frac{1}{2} \|\mathbf{w}\|^2$$

is similar to that based on logistic regression.

- The corresponding regularized loss takes the form:

$$\frac{1}{N} \sum_{i=1}^N (1 - y_i [w_0 + \mathbf{x}_i^T \mathbf{w}])^+ + \frac{1}{NC} \|\mathbf{w}\|^2,$$

- Here the regularization parameter is $\lambda = \frac{1}{NC}$

- The first term above corresponds to the negative log likelihood:

$$\frac{1}{N} \sum_{i=1}^N -\log g(y_i [w_0 + \mathbf{x}_i^T \mathbf{w}]) \quad , g(z) = \frac{1}{1 + e^{-z}} = \text{logistic function}$$

Probabilistic Interpretation of SVM

- The corresponding regularized loss takes the form:

$$\frac{1}{N} \sum_{i=1}^N (1 - y_i [w_0 + \mathbf{x}_i^T \mathbf{w}])^+ + \frac{1}{NC} \|\mathbf{w}\|^2,$$

- One in the pursue of a probabilistic interpretation of SVM can show:

$$\begin{aligned} & \exp(-2(1 - y_i [w_0 + \mathbf{x}_i^T \mathbf{w}])^+) \\ &= \int_0^\infty \frac{1}{\sqrt{2\pi\lambda_i}} \exp\left(-\frac{1}{2} \frac{(1 + \lambda_i - y_i [w_0 + \mathbf{x}_i^T \mathbf{w}])^2}{\lambda_i}\right) d\lambda_i \end{aligned}$$

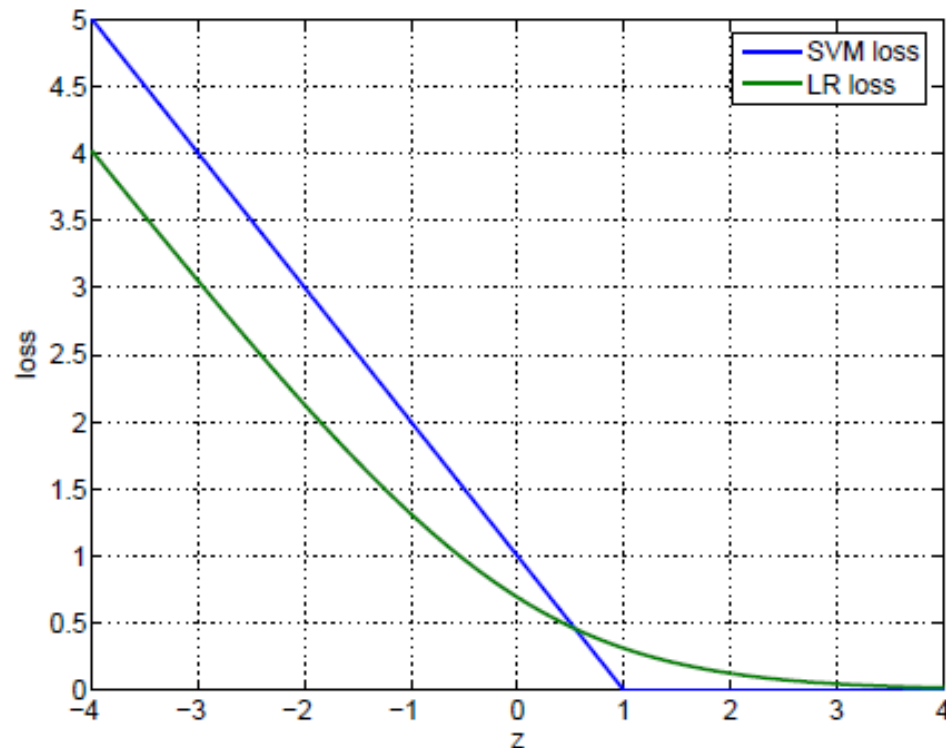
- This allows to fit the SVM with EM using λ_i as latent variables.

- Sollich, P. (2002). [Bayesian methods for support vector machines: evidence and predictive class probabilities](#). *Machine Learning* 46, 21– 52.
- Polson, N. and S. Scott (2011). [Data augmentation for support vector machines](#). *Bayesian Analysis* 6(1), 1–124.
- Franc, V., A. Zien, and B. Schoelkopf (2011). [Support vector machines as probabilistic models](#). In *Intl. Conf. on Machine Learning*.

SVM Vs. Logistic Regression

□ Common objective the minimization of:

$$\frac{1}{N} \sum_{i=1}^N \text{Loss}(y_i [w_0 + \mathbf{x}_i^T \mathbf{w}]) + \frac{1}{NC} \|\mathbf{w}\|^2,$$



$$\text{SVM Loss: } \text{Loss}(z) = (1 - z)^+$$

$$\text{LR Loss: } \text{Loss}(z) = \frac{1}{1 + e^{-z}}$$

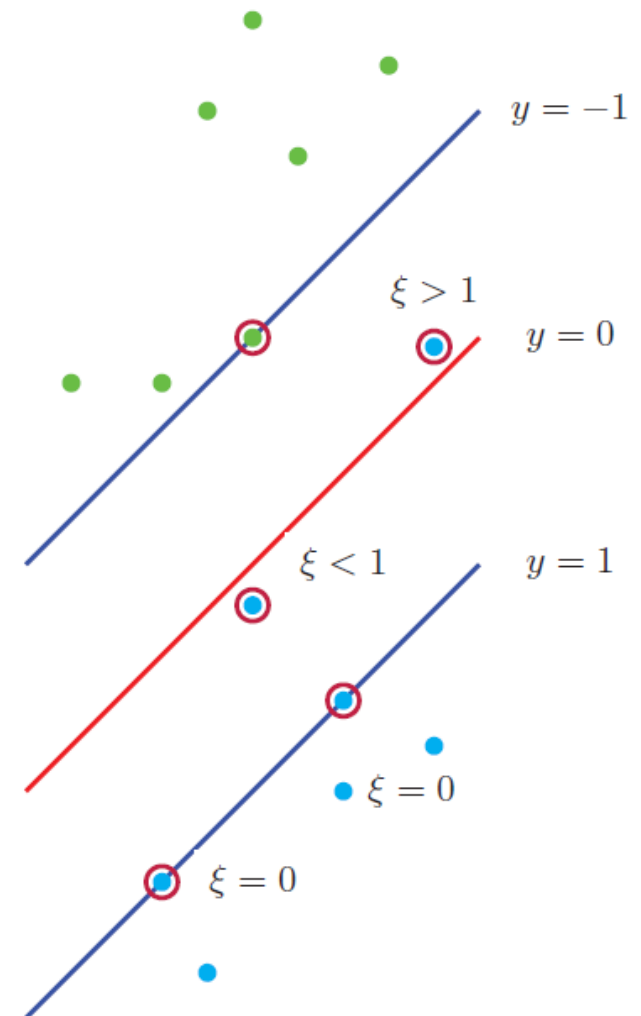
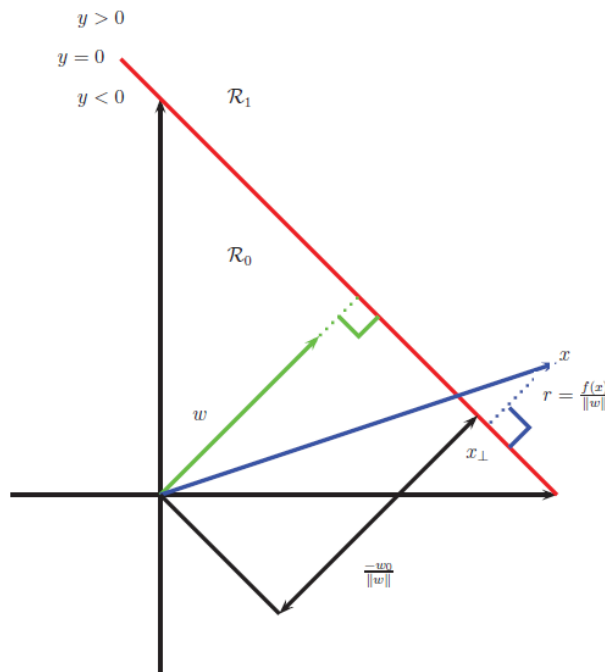
SVM - Separable Case

□ Our problem is stated as follows:

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 \text{ subject to}$$

$$y_i [w_0 + \mathbf{x}_i^T \mathbf{w}] - 1 \geq 0, i = 1, \dots, N$$

(each point is on the correct side of the boundary)



SVM - Separable Case

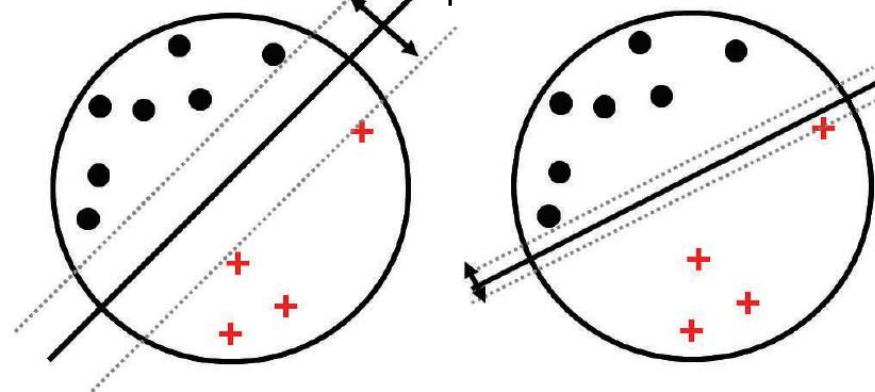
- Our problem is stated as follows:

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 \text{ subject to}$$

$$y_i [w_0 + \mathbf{x}_i^T \mathbf{w}] - 1 \geq 0, i = 1, \dots, N$$

(each point is on the correct side of the boundary)

Margin here defined as the perpendicular distance to the closest point



Large margin

Small margin

- We represent the constraints with Lagrange multipliers as losses:

$$\max_{\alpha \geq 0} \alpha_i (1 - y_i [w_0 + \mathbf{x}_i^T \mathbf{w}]) = \begin{cases} 0, & \text{if } y_i [w_0 + \mathbf{x}_i^T \mathbf{w}] - 1 \geq 0 \\ \infty, & \text{otherwise} \end{cases}$$

- The minimization problem can now be stated as:

$$\min_{\mathbf{w}} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^N \max_{\alpha \geq 0} \alpha_i (1 - y_i [w_0 + \mathbf{x}_i^T \mathbf{w}]) \right\}$$

SVM - Separable Case

- The minimization problem can now be stated as:

$$\min_{\mathbf{w}} \max_{\{\alpha_i \geq 0\}} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^N \alpha_i (1 - y_i [w_0 + \mathbf{x}_i^T \mathbf{w}]) \right\} =$$

$$\max_{\{\alpha_i \geq 0\}} \min_{\mathbf{w}} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^N \alpha_i (1 - y_i [w_0 + \mathbf{x}_i^T \mathbf{w}]) \right\}$$

- As a result we should be able to minimize wrt \mathbf{w} for any given set of the α 's the following:

$$J(\mathbf{w}; \alpha) = \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^N \alpha_i (1 - y_i [w_0 + \mathbf{x}_i^T \mathbf{w}]) \right\}$$

- Setting the derivatives wrt to \mathbf{w} and w_0 gives:

$$\mathbf{w} - \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i = 0, \quad -\sum_{i=1}^N \alpha_i y_i = 0$$

SVM - Separable Case

- Substituting these Eqs $\hat{\mathbf{w}} - \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i = 0$, $-\sum_{i=1}^N \alpha_i y_i = 0$ into the objective function

$$\max_{\{\alpha_i \geq 0\}} \min_{\mathbf{w}} \left\{ \frac{1}{2} \|\hat{\mathbf{w}}\|^2 + \sum_{i=1}^N \alpha_i (1 - y_i [\mathbf{w}_0 + \mathbf{x}_i^T \hat{\mathbf{w}}]) \right\}$$

gives:

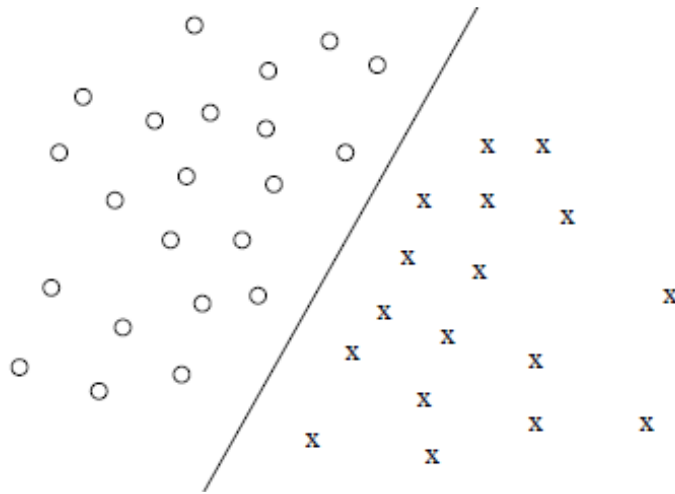
$$\max_{\{\alpha_i \geq 0, \sum_{i=1}^N \alpha_i y_i = 0\}} \left\{ \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N y_i y_j \alpha_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j \right\}$$

- Only α_i 's corresponding to “support vectors” will be non-zero.
- Once this optimization problem is solved, can do predictions on any test \mathbf{x} according to the sign of the discriminant function

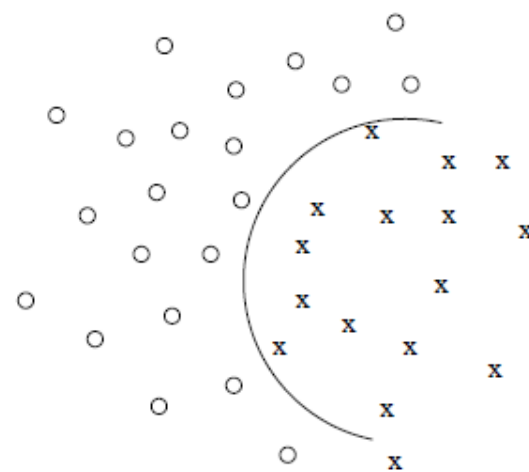
$$\hat{\mathbf{w}}_0 + \mathbf{x}^T \mathbf{w} = \hat{\mathbf{w}}_0 + \mathbf{x}^T \sum_{i,j=1}^N \hat{\alpha}_i y_i \mathbf{x}_i = \hat{\mathbf{w}}_0 + \sum_{i \in SV} \hat{\alpha}_i y_i (\mathbf{x}^T \mathbf{x}_i)$$

Non-linear Classifier

- So far our classifier can make only linear separations.
- We can easily obtain a non-linear classifier by using the kernel trick.
Map $\mathbf{x} = [x_1 \ x_2]$ into $\boldsymbol{\phi}(\mathbf{x}) = [x_1^2 \ x_2^2 \ \sqrt{2}x_1x_2 \ \sqrt{2}x_1 \ \sqrt{2}x_2 \ 1]$



Linear separator in the
feature $\boldsymbol{\phi}$ -space

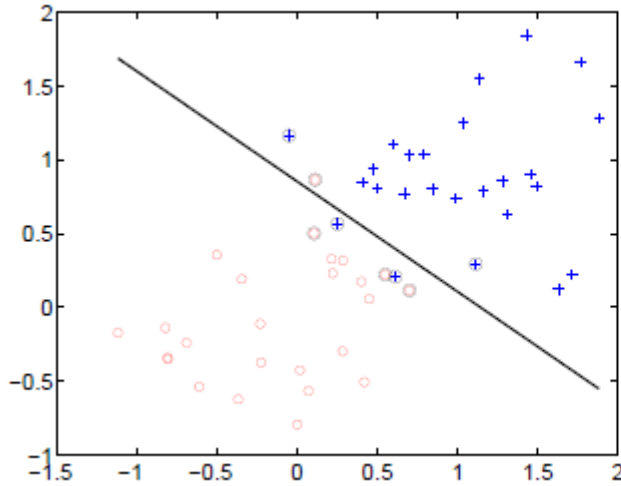


Non-linear separator
in the original \mathbf{x} -space

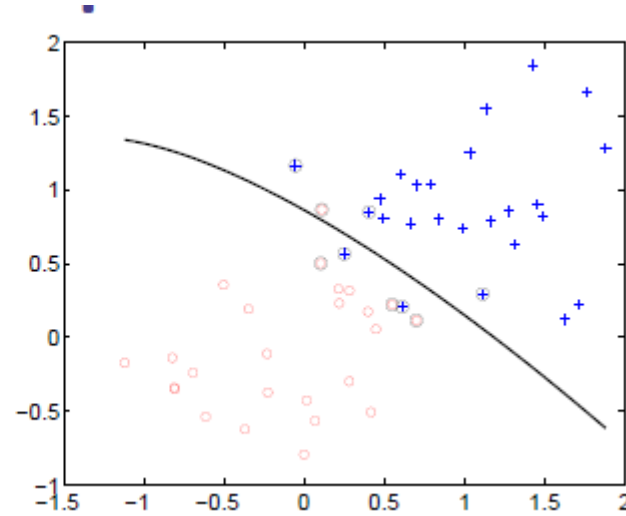
- The algorithm only needs dot products of the form:
$$\boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{x}') = x_1^2 x_1'^2 + x_2^2 x_2'^2 + 2x_1 x_2 x_1' x_2' + 2x_1 x_1' + 2x_2 x_2' + 1 = (1 + \mathbf{x}^T \mathbf{x}')^2$$
- So the inner products can be evaluated without ever explicitly constructing the feature vectors $\boldsymbol{\phi}(\mathbf{x})$!

SVM

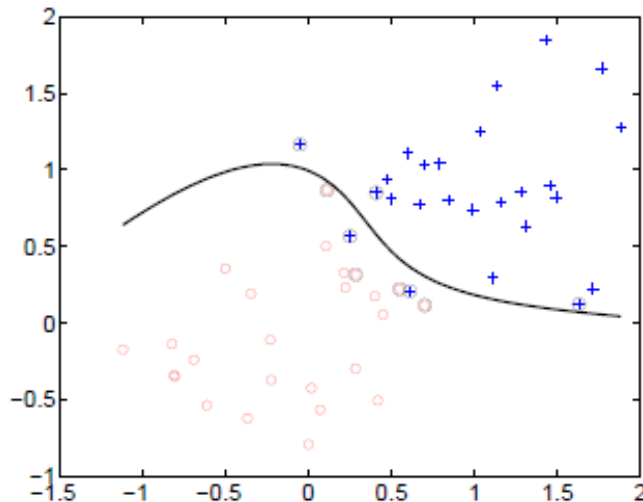
linear



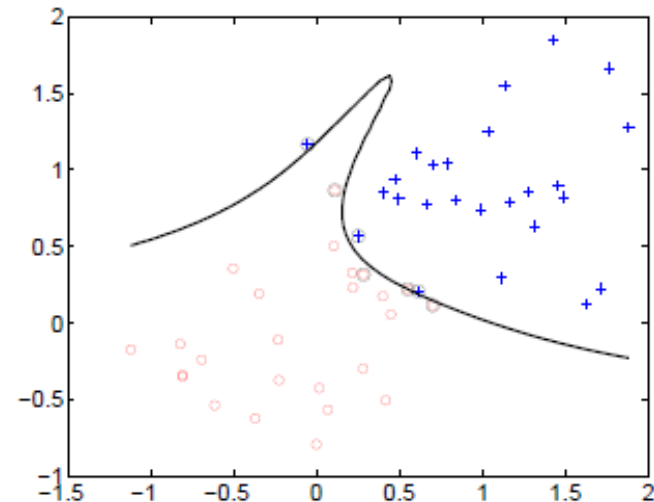
2nd order polynomial



4th order polynomial

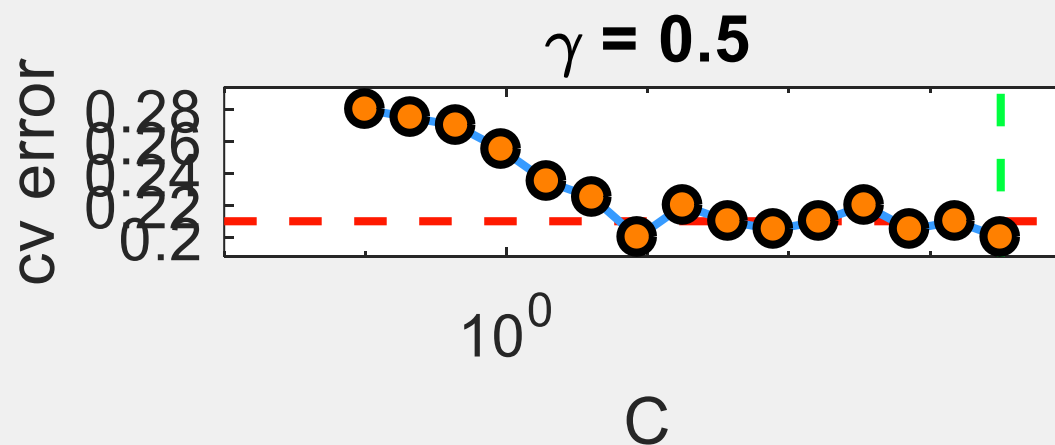
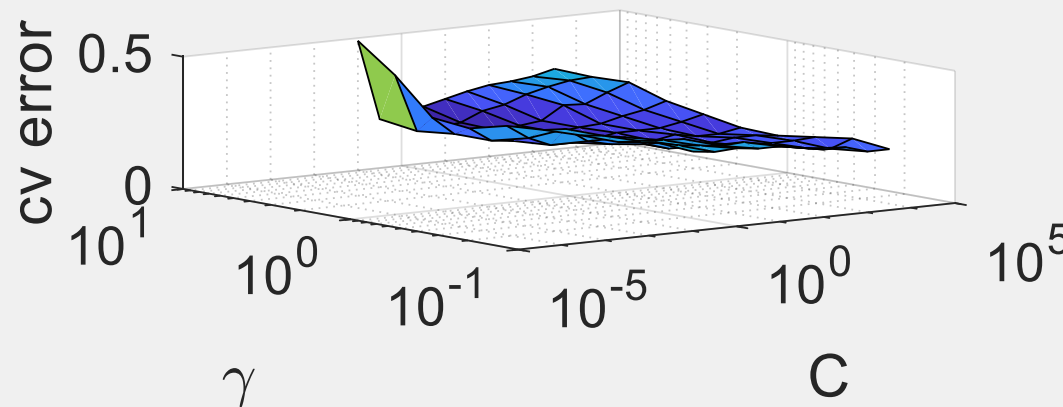


8th order polynomial



Choosing C

- ❑ Need to specify the kernel parameters and C . Typically C is chosen by cross-validation.
- ❑ C interacts quite strongly with the kernel parameters.
- ❑ Consider an RBF kernel with precision $\gamma = \frac{1}{(2\sigma)^2}$
- ❑ If $\gamma = 5$, corresponding to narrow kernels, we need heavy regularization, and hence small C (so $\lambda = 1/C$ is big).
- ❑ If $\gamma = 1$, a larger value of C should be used.



CV estimate of the 0-1 risk as a function of C and γ .

[svmCGammaDemo.m](#)
from [PMTK3](#)

Choosing C

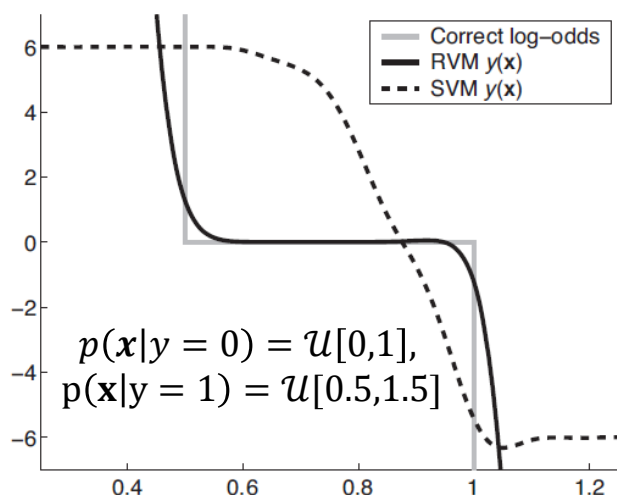
- ❑ libsvm recommends using CV over a 2d grid with values $C \in \{2^{-5}, 2^{-3}, \dots, 2^{15}\}$ and $\gamma \in \{2^{-15}, 2^{-13}, \dots, 2^3\}$.
- ❑ Standardize the data first, for a spherical Gaussian kernel to make sense.
- ❑ To choose C efficiently, one can develop a path following the spirit of lars.
 - The basic idea is to start with λ large, so that the margin $1/||w(\lambda)||$ is wide, and hence all points are inside of it and have $\alpha_i = 1$.
 - By slowly decreasing λ , a small set of points will move from inside the margin to outside, and their α_i values will change from 1 to 0, as they cease to be support vectors.
 - When λ is maximal, the function is completely smoothed, and no support vectors remain.
- Hastie, T., S. Rosset, R. Tibshirani, and J. Zhu (2004). [The entire regularization path for the support vector machine](#). *J. of Machine Learning Research* 5, 1391–1415.

Probabilistic Response

- ❑ An SVM classifier produces a hard-labeling, $\hat{y}(x) = \text{sign}(f(x))$.
- ❑ A heuristic approach to producing confidence in our prediction is to interpret $f(x)$ as the log-odds ratio, $\log \frac{p(y = 1|x)}{p(y = 0|x)}$ and convert the output of an SVM to a probability using

$$p(y = 1|x, \theta) = \sigma(af(x) + b)$$

where a, b can be estimated by maximum likelihood on a separate validation set (using the training set to estimate a and b leads to severe overfitting.)



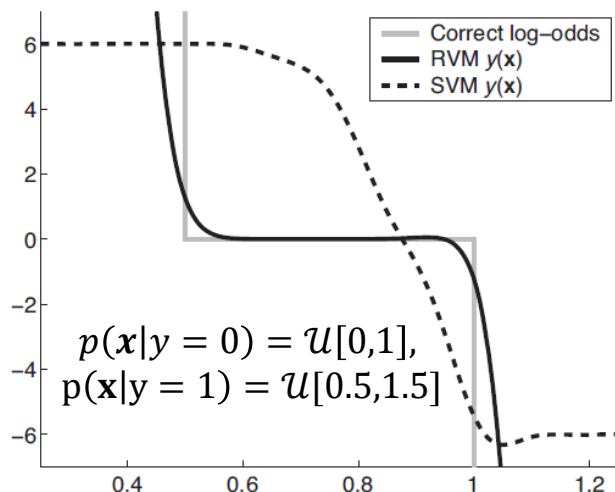
- The log-odds of class 1 over class 0 should be zero in $[0.5, 1.0]$, and infinite outside this region.
- We sampled 1000 points from the model, and then fit an RVM and an SVM with a Gaussian kernel of width 0.1.
- Both models can perfectly capture the decision boundary, and achieve a generalization error of 25%, which is Bayes optimal in this problem.
- The probabilistic output from the RVM is a good approximation to the true log-odds, but this is not the case for the SVM.

Probabilistic Response

- An SVM classifier produces a hard-labeling, $\hat{y}(\mathbf{x}) = \text{sign}(f(\mathbf{x}))$.
- A heuristic approach to producing confidence in our prediction is to interpret $f(\mathbf{x})$ as the log-odds ratio, $\log \frac{p(y = 1|\mathbf{x})}{p(y = 0|\mathbf{x})}$ and convert the output of an SVM to a probability using

$$p(y = 1|\mathbf{x}, \boldsymbol{\theta}) = \sigma(af(\mathbf{x}) + b)$$

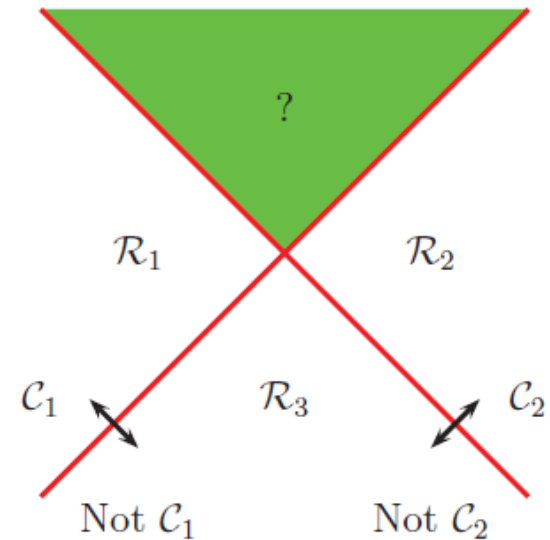
where a , b can be estimated by maximum likelihood on a separate validation set (using the training set to estimate a and b leads to severe overfitting.)



- Tipping, M. (2001). [Sparse Bayesian learning and the relevance vector machine](#). *J. of Machine Learning Research* 1, 211–244.

Multiclass Classification

- ❑ The obvious approach is to use a **one-vs-all** approach. We train C binary classifiers, $f_c(x)$, where the data from class c is treated as positive, and the data from all the other classes as negative.
- ❑ This can result in ambiguously labeled regions.
- ❑ Can pick $\hat{y}(x) = \underset{c}{\operatorname{argmax}} f_c(x)$. This will not work since the f_c functions don't have comparable magnitudes.

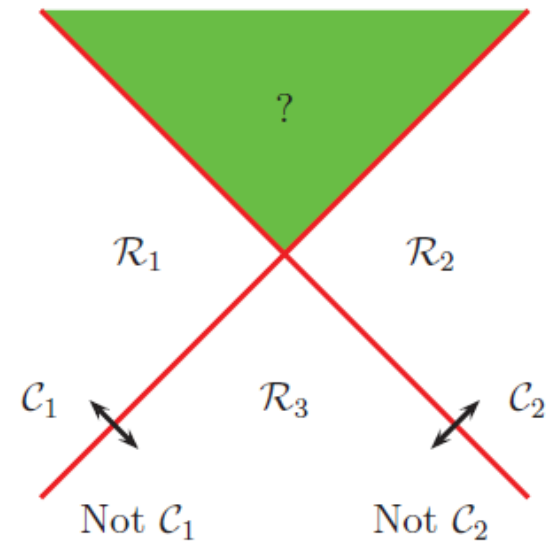


The one-versus-rest approach.
The green region is predicted to be both class 1 and class 2.

- Weston, J. and C. Watkins (1999). [Multi-class support vector machines](#). In *ESANN*.

Multiclass Classification

- ❑ Each binary subproblem is likely to suffer from the **class imbalance** problem.
- ❑ E.g. consider 10 equally represented classes. When training f_1 , we will have 10% positive examples and 90% negative examples, which can hurt performance.
- ❑ One can train all C classifiers simultaneously, but the resulting method takes $\mathcal{O}(C^2 N^2)$ time, instead of $\mathcal{O}(CN^2)$.

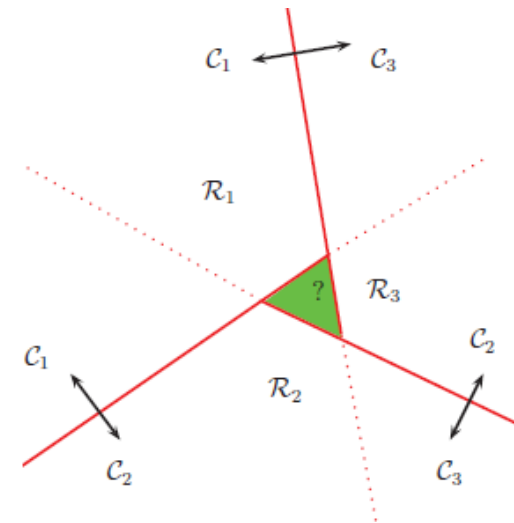


The one-versus-rest approach.
The green region is predicted to be both class 1 and class 2.

- Weston, J. and C. Watkins (1999). [Multi-class support vector machines](#). In *ESANN*.

Multiclass Classification

- ❑ Another approach is the **one-versus-one**.
- ❑ We train $C(C - 1)/2$ classifiers to discriminate all pairs $f_{c,c'}$. We then classify a point into the class which has the highest number of votes.
- ❑ This can also result in ambiguities.
- ❑ It takes $\mathcal{O}(C^2 N^2)$ time to train and $\mathcal{O}(C^2 N_{sv})$ to test each data point, where N_{sv} is the number of support vectors.²
- ❑ Fundamental limitation of SVMs: they do not model uncertainty, so their output scores are not comparable across classes..



The one-versus-one approach.

- Allwein, E., R. Schapire, and Y. Singer (2000). [Reducing multiclass to binary: A unifying approach for margin classifiers](#). *J. of Machine Learning Research*, 113–141.

Relevance Vector Machines

Relevance Vector Machines

- ❑ Relevance Vector Machines (RVM) were introduced by [M. Tipping](#) in [Tipping, M. E. \(2001\). Sparse Bayesian learning and the relevance vector machine. Journal of Machine Learning Research 1, 211–244.](#)
- ❑ It is a *Bayesian regression* technique on a *generalized linear model* that produces *sparse representations* (in the sense that many of the coefficients of the generalized linear model turn out to be zero).
- ❑ We will first describe it on one output dimension and then generalize it to many dimensions.
- ❑ Finally, we will couple it with an adaptive input decomposition algorithm and apply it to uncertainty quantification tasks. [See I. Bilonis and N. Zabarar. Multidimensional adaptive relevance vector machines for uncertainty quantification. SIAM Journal for Scientific Computing, Vol. 34, No. 6, pp. B881–B908 2012](#) for further details.

Relevance Vector Machines

- ❑ RVM addresses many of the limitations of SVM:
 - SVM outputs are decisions rather than posterior probabilities.
 - Extension to $K > 2$ classes in SVM is problematic.
 - There are complexity parameters in SVM to be found by cross-validation.
 - Predictions in SVM are based on kernel functions centered at the training data that are required to be positive definite.

RVM: Likelihood

- We assume that we observe the data:

$$\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)} = f(\mathbf{x}^{(i)})\}_{i=1}^n,$$

and that we wish to learn the function $f(\cdot)$.

- Pick a set of M basis functions $\boldsymbol{\phi}: \mathbf{R}^k \rightarrow \mathbf{R}^M$:

$$\boldsymbol{\phi}(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_M(\mathbf{x})).$$

- We will approximate $f(\cdot)$ by:

$$\hat{f}(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) + \epsilon = \sum_{j=1}^M w_j \phi_j(\mathbf{x}) + \epsilon,$$

where ϵ is a random variable representing noise.

- For simplicity, we choose to work with Gaussian noise with zero mean and inverse variance β , so the **likelihood** is:

$$p(t|\mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t|\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}), \beta^{-1})$$

RVM: Prior

- The essence of RVM is the following **prior on the weights**:

$$p(\mathbf{w}|\boldsymbol{\alpha}) = \prod_{j=1}^M p(w_j|\alpha_j),$$

with $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_M)$ and:

$$p(w_j|\alpha_j) = \mathcal{N}(w_j|0, \alpha_j^{-1}).$$

- The characteristic of this prior is that:

$$\lim_{\alpha_j \rightarrow \infty} p(w_j|\alpha_j) = \delta(w_j).$$

- This means that the j -th basis function can be removed from the model if its α_j is very big. This is because, its weights will be very sharply picked about zero.
- To complete the model, we also need to assign a prior $p(\boldsymbol{\beta}, \boldsymbol{\alpha})$. We manage to get rid of it by employing the *evidence approximation*.

RVM: Posterior

□ The **posterior of the weights** is:

$$p(\mathbf{w}|\mathcal{D}, \beta, \alpha) = \mathcal{N}(\mathbf{w}|\mathbf{m}, \Sigma),$$

where:

$$\Sigma = (\text{diag}(\alpha) + \beta \Phi^T \Phi)^{-1},$$
$$\mathbf{m} = \beta \Sigma \Phi^T \mathbf{t},$$

with $\mathbf{t} = (t^{(1)}, \dots, t^{(n)})$ is the *vector of observations*, $\Phi \in \mathbf{R}^{N \times M}$ is the *design matrix*:

$$\Phi_{ij} = \phi_j(\mathbf{x}^{(i)}),$$

and $\text{diag}(\alpha)$ is a diagonal matrix with α on the diagonal.

Proof:

First, notice that we can write:

$$p(\mathbf{w}|\alpha) = \prod_{j=1}^M p(w_j|\alpha_j) = \prod_{j=1}^M \mathcal{N}(w_j|0, \alpha_j^{-1}) = \mathcal{N}_M(\mathbf{w}|\text{diag}(\alpha)^{-1}).$$

RVM: Posterior

Proof (continuation):

Also, notice that the likelihood of the data (under an independence assumption) is given by:

$$p(\mathcal{D}|\mathbf{w}, \beta) = \prod_{i=1}^n \mathcal{N}(t^{(i)} | \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}^{(i)}), \beta^{-1}) = \mathcal{N}_N(\mathbf{y} | \boldsymbol{\Phi} \mathbf{w}, \beta^{-1} \mathbf{I}_N),$$

where \mathbf{I}_N is the N -dimensional unit matrix and we have used the fact that:

$$\begin{aligned} \sum_{i=1}^N \| \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}^{(i)}) - t^{(i)} \|^2 &= \sum_{i=1}^N (\boldsymbol{\phi}(\mathbf{x}^{(i)}) \mathbf{w} - t^{(i)})^T (\boldsymbol{\phi}(\mathbf{x}^{(i)}) \mathbf{w} - t^{(i)}) \\ &= (\boldsymbol{\Phi} \mathbf{w} - \mathbf{y})^T (\boldsymbol{\Phi} \mathbf{w} - \mathbf{y}). \end{aligned}$$

Finally, we have:

$$p(\mathbf{w} | \mathcal{D}, \boldsymbol{\alpha}, \beta) \propto p(\mathcal{D} | \mathbf{w}, \beta) p(\mathbf{w} | \boldsymbol{\alpha}) = \mathcal{N}_N(\mathbf{y} | \boldsymbol{\Phi} \mathbf{w}, \beta^{-1} \mathbf{I}_N) \mathcal{N}_M(\mathbf{w} | \text{diag}(\boldsymbol{\alpha})^{-1}).$$

RVM: Posterior

Proof (continuation):

Remember that the posterior is:

$$p(\mathbf{w}|\mathcal{D}, \boldsymbol{\alpha}, \beta) \propto p(\mathcal{D}|\mathbf{w}, \beta)p(\mathbf{w}|\boldsymbol{\alpha}) = \mathcal{N}_N(\mathbf{y}|\boldsymbol{\Phi}\mathbf{w}, \beta^{-1}\mathbf{I}_N)\mathcal{N}_M(\mathbf{w}|\text{diag}(\boldsymbol{\alpha})^{-1}).$$

Now, we just have to look at it and complete the square...

You can simply drop this guy!

We have for the exponential inside the right hand side:

$$\begin{aligned} & (\boldsymbol{\Phi}\mathbf{w} - \mathbf{t})^T(\beta\mathbf{I}_N)(\boldsymbol{\Phi}\mathbf{w} - \mathbf{t}) + \mathbf{w}^T\text{diag}(\boldsymbol{\alpha})\mathbf{w} = \\ & \mathbf{w}^T(\beta\boldsymbol{\Phi}^T\boldsymbol{\Phi})\mathbf{w} - 2\mathbf{w}^T(\beta\boldsymbol{\Phi}^T)\mathbf{y} + \beta\mathbf{t}^T\mathbf{t} + \mathbf{w}^T\text{diag}(\boldsymbol{\alpha})\mathbf{w} = \\ & \mathbf{w}^T(\text{diag}(\boldsymbol{\alpha}) + \beta\boldsymbol{\Phi}^T\boldsymbol{\Phi})\mathbf{w} - 2\mathbf{w}^T(\beta\boldsymbol{\Phi}^T)\mathbf{t} = \\ & \mathbf{w}^T\boldsymbol{\Sigma}^{-1}\mathbf{w} - 2\mathbf{w}^T\boldsymbol{\Sigma}^{-1}(\beta\boldsymbol{\Sigma}\boldsymbol{\Phi}^T)\mathbf{t}. \end{aligned}$$

From which the result follows directly. End of Proof.

Important remark: When completing the square for a distribution, you can subtract, add or even ignore anything that does not depend on the random variable whose posterior you are calculating!

Or using Results for Linear Gaussian Models

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$$

$$p(\mathbf{y} \mid \mathbf{x}) = \mathcal{N}(\mathbf{y} \mid \mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1})$$

□ For the above linear Gaussian model, the following hold:

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} \mid \mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T)$$

$$p(\mathbf{x} \mid \mathbf{y}) = \mathcal{N}\left(\mathbf{x} \mid \left(\boldsymbol{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A}\right)^{-1} \left(\boldsymbol{\Lambda} \boldsymbol{\mu} + \mathbf{A}^T \mathbf{L}(\mathbf{y} - \mathbf{b})\right), \left(\boldsymbol{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A}\right)^{-1}\right)$$

□ For our problem $\mathbf{x} \rightarrow \mathbf{w}$, $\boldsymbol{\mu} \rightarrow \mathbf{0}$, $\boldsymbol{\Lambda} \rightarrow \text{diag}(\boldsymbol{\alpha})$, $\mathbf{y} \rightarrow \mathbf{t}$, $\mathbf{A} \rightarrow \boldsymbol{\Phi}$, $\mathbf{b} \rightarrow \mathbf{0}$, $\mathbf{L} = \beta \mathbf{I}_N$

$$p(\mathbf{t} \mid \mathbf{w}) = \mathcal{N}_N(\mathbf{y} \mid \boldsymbol{\Phi} \mathbf{w}, \beta^{-1} \mathbf{I}_N), p(\mathbf{w}) = \mathcal{N}_M(\mathbf{w} \mid \text{diag}(\boldsymbol{\alpha})^{-1})$$

$$\text{and: } p(\mathbf{w} \mid \mathcal{D}, \beta, \boldsymbol{\alpha}) = \mathcal{N}\left(\mathbf{w} \mid \beta \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \mathbf{t}, \underbrace{\left(\text{diag}(\boldsymbol{\alpha}) + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}\right)^{-1}}_{\boldsymbol{\Sigma}}\right)$$

RVM: Predictive distribution

The **predictive** distribution is easily found as follows:

$$\begin{aligned} p(t|\mathbf{x}, \mathbf{X}, \mathbf{t}, \boldsymbol{\alpha}, \beta) &= \int p(t|\mathbf{x}, \mathbf{w}, \beta) p(\mathbf{w}|\mathbf{X}, \mathbf{t}, \boldsymbol{\alpha}, \beta) d\mathbf{w} \\ &= \mathcal{N}(t|\mathbf{m}^T \boldsymbol{\phi}(\mathbf{x}), \boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\Sigma} \boldsymbol{\phi}(\mathbf{x}) + \beta^{-1}). \end{aligned}$$

So, we have the *predictive mean* (\mathbf{m} being the posterior mean):

$$\mu(\mathbf{x}) = \mathbf{m}^T \boldsymbol{\phi}(\mathbf{x}),$$

and the predictive variance:

$$\sigma^2(\mathbf{x}) = \boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\Sigma} \boldsymbol{\phi}(\mathbf{x}) + \beta^{-1}.$$

Proof:

Complete the square....

End of proof.

RVM: Evidence Approximation

- We will select the hyper-parameters α and β by maximizing employing the *evidence approximation* (otherwise known as *maximizing the marginal likelihood*). We start by motivating it.
- Assume that we have a prior for α and β , $p(\alpha, \beta)$. Then the posterior of these parameters is:

$$p(\alpha, \beta | \mathcal{D}) \propto p(\mathcal{D} | \alpha, \beta) p(\alpha, \beta) = \int p(\mathcal{D} | \mathbf{w}, \beta) p(\mathbf{w} | \alpha) d\mathbf{w} p(\alpha, \beta).$$



Marginal Likelihood (or Evidence)

- Under the evidence approximation, we assume that:
 - $p(\alpha, \beta)$ is relatively flat.
 - $p(\mathcal{D} | \alpha, \beta)$ has very sharp maxima.

Then, we may obtain point estimates of the hyper-parameters by identifying these maxima.

RVM: Evidence

- To carry out the evidence approximation, we must obtain an analytic expression for the marginal likelihood. This is possible for RVM:

$$\mathcal{E}(\boldsymbol{\alpha}, \beta) = \log(p(\mathcal{D}|\boldsymbol{\alpha}, \beta)) = -\frac{1}{2}[N \log(2\pi) + \log|\mathbf{C}| + \mathbf{t}^T \mathbf{C}^{-1} \mathbf{t}],$$

where $\mathbf{C} \in \mathbf{R}^{N \times N}$, is given by:

$$\mathbf{C} = \beta^{-1} \mathbf{I}_N + \boldsymbol{\Phi} \text{diag}(\boldsymbol{\alpha})^{-1} \boldsymbol{\Phi}^T.$$

Proof:

This does require some more work that we avoided doing before. We start again from:

$$p(\mathcal{D}|\mathbf{w}, \beta)p(\mathbf{w}|\boldsymbol{\alpha}) = \mathcal{N}_N(\mathbf{t}|\boldsymbol{\Phi}\mathbf{w}, \beta^{-1}\mathbf{I}_N)\mathcal{N}_M(\mathbf{w}|\text{diag}(\boldsymbol{\alpha})^{-1}),$$

and we will complete the square making sure this time that we keep track of all terms that depend on \mathbf{t} .

RVM: Evidence

Proof (continuation):

Remember, that we want to integrate \mathbf{w} out of:

$$p(\mathcal{D}|\mathbf{w}, \beta)p(\mathbf{w}|\boldsymbol{\alpha}) = \mathcal{N}_n(\mathbf{t}|\boldsymbol{\Phi}\mathbf{w}, \beta^{-1}\mathbf{I}_N)\mathcal{N}_M(\mathbf{w}|\text{diag}(\boldsymbol{\alpha})^{-1}).$$

Again, we look at whatever is inside the exponential of the right-hand side and attempt to complete the square. This time, we already know the answer, but we need to keep track of \mathbf{t} explicitly:

$$\begin{aligned} & (\boldsymbol{\Phi}\mathbf{w} - \mathbf{t})^T(\beta\mathbf{I}_N)(\boldsymbol{\Phi}\mathbf{w} - \mathbf{t}) + \mathbf{w}^T\text{diag}(\boldsymbol{\alpha})\mathbf{w} = \\ & \mathbf{w}^T(\beta\boldsymbol{\Phi}^T\boldsymbol{\Phi})\mathbf{w} - 2\mathbf{w}^T(\beta\boldsymbol{\Phi}^T)\mathbf{t} + \beta\mathbf{t}^T\mathbf{t} + \mathbf{w}^T\text{diag}(\boldsymbol{\alpha})\mathbf{w} = \\ & \mathbf{w}^T(\text{diag}(\boldsymbol{\alpha}) + \beta\boldsymbol{\Phi}^T\boldsymbol{\Phi})\mathbf{w} - 2\mathbf{w}^T(\beta\boldsymbol{\Phi}^T)\mathbf{t} + \beta\mathbf{t}^T\mathbf{t} = \\ & \mathbf{w}^T\boldsymbol{\Sigma}^{-1}\mathbf{w} - 2\mathbf{w}^T\boldsymbol{\Sigma}^{-1}(\beta\boldsymbol{\Sigma}\boldsymbol{\Phi}^T)\mathbf{t} + \beta\mathbf{t}^T\mathbf{t} = \\ & \mathbf{w}^T\boldsymbol{\Sigma}^{-1}\mathbf{w} - 2\mathbf{w}^T\boldsymbol{\Sigma}^{-1}\mathbf{m} + \mathbf{m}^T\boldsymbol{\Sigma}^{-1}\mathbf{m} - \mathbf{m}^T\boldsymbol{\Sigma}^{-1}\mathbf{m} + \beta\mathbf{t}^T\mathbf{t} = \\ & (\mathbf{w} - \mathbf{m})^T\boldsymbol{\Sigma}^{-1}(\mathbf{w} - \mathbf{m}) + (\beta\mathbf{t}^T\mathbf{t} - \mathbf{m}^T\boldsymbol{\Sigma}^{-1}\mathbf{m}). \end{aligned}$$

RVM: Evidence

Proof (continuation):

We can now write:

$$p(\mathcal{D}|\mathbf{w}, \beta)p(\mathbf{w}|\boldsymbol{\alpha}) = \mathcal{N}_N(\mathbf{y}|\boldsymbol{\Phi}\mathbf{w}, \beta^{-1}\mathbf{I}_N)\mathcal{N}_M(\mathbf{w}|\text{diag}(\boldsymbol{\alpha})^{-1}) = \\ (2\pi)^{-\frac{N}{2}}\beta^{\frac{N}{2}}(2\pi)^{-\frac{M}{2}}|\text{diag}(\boldsymbol{\alpha})|^{\frac{1}{2}} \exp\left\{-\frac{1}{2}(\mathbf{w} - \mathbf{m})^T \boldsymbol{\Sigma}^{-1}(\mathbf{w} - \mathbf{m})\right\}$$

RVM: Two useful matrix identities

□ Woodbury Matrix Identity:

$$(A + UCV)^{-1} = A^{-1} - A^{-1}U(C^{-1} + VA^{-1}U)^{-1}VA^{-1}.$$

□ Matrix Determinant Lemma:

$$|A + UWV^T| = |W^{-1} + V^T A^{-1}U| |W| |A|.$$

You do not have to remember these, but you must be aware of their existence. Sometimes, they are the only route to analytic progress!

Or using Results for Linear Gaussian Models

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$$

$$p(\mathbf{y} \mid \mathbf{x}) = \mathcal{N}(\mathbf{y} \mid \mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1})$$

□ For the above linear Gaussian model, the following hold:

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} \mid \mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T)$$

$$p(\mathbf{x} \mid \mathbf{y}) = \mathcal{N}\left(\mathbf{x} \mid \left(\boldsymbol{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A}\right)^{-1} \left(\boldsymbol{\Lambda} \boldsymbol{\mu} + \mathbf{A}^T \mathbf{L}(\mathbf{y} - \mathbf{b})\right), \left(\boldsymbol{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A}\right)^{-1}\right)$$

□ For our problem $\mathbf{x} \rightarrow \mathbf{w}$, $\boldsymbol{\mu} \rightarrow \mathbf{0}$, $\boldsymbol{\Lambda} \rightarrow \text{diag}(\boldsymbol{\alpha})$, $\mathbf{y} \rightarrow \mathbf{t}$, $\mathbf{A} \rightarrow \boldsymbol{\Phi}$, $\mathbf{b} \rightarrow \mathbf{0}$, $\mathbf{L} = \beta \mathbf{I}_N$

$$p(\mathbf{t} \mid \mathbf{w}) = \mathcal{N}_N(\mathbf{y} \mid \boldsymbol{\Phi} \mathbf{w}, \beta^{-1} \mathbf{I}_N), p(\mathbf{w}) = \mathcal{N}_M(\mathbf{w} \mid \text{diag}(\boldsymbol{\alpha})^{-1})$$

$$\text{and: } p(\mathbf{t} \mid \beta, \boldsymbol{\alpha}) = \mathcal{N}(\mathbf{w} \mid \mathbf{0}, \beta^{-1} \mathbf{I}_N + \boldsymbol{\Phi} \text{diag}(\boldsymbol{\alpha})^{-1} \boldsymbol{\Phi}^T)$$

RVM: Evidence Approximation

□ From [our earlier derivation](#):

$$p(\mathcal{D}|\boldsymbol{\alpha}, \beta) = (2\pi)^{-\frac{N}{2}} \beta^{\frac{N}{2}} |\text{diag}(\boldsymbol{\alpha})|^{\frac{1}{2}} \exp \left\{ -\frac{1}{2} \mathbf{t}^T (\beta \mathbf{I}_N - \beta^2 \boldsymbol{\Phi} \boldsymbol{\Sigma} \boldsymbol{\Phi}^T) \mathbf{t} \right\} |\boldsymbol{\Sigma}|^{\frac{1}{2}}$$

$$\text{where } \boldsymbol{\Sigma} = (\text{diag}(\boldsymbol{\alpha}) + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1}, \quad \mathbf{m} = \beta \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \mathbf{t}$$

□ Noting that $\mathbf{t}^T \beta^2 \boldsymbol{\Phi} \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \mathbf{t} = \mathbf{m}^T \boldsymbol{\Sigma}^{-1} \mathbf{m}$, can simplify as:

$$p(\mathcal{D}|\boldsymbol{\alpha}, \beta) = (2\pi)^{-\frac{N}{2}} \beta^{\frac{N}{2}} |\text{diag}(\boldsymbol{\alpha})|^{\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\beta \mathbf{t}^T \mathbf{t} - \mathbf{m}^T \boldsymbol{\Sigma}^{-1} \mathbf{m}) \right\} |\boldsymbol{\Sigma}|^{\frac{1}{2}}$$

□ Taking derivatives wrt α_i of $\ln p(\mathcal{D}|\boldsymbol{\alpha}, \beta)$ and using $\frac{\partial \ln |\boldsymbol{\Sigma}|}{\partial \alpha_i} =$

$$\text{Tr} \left(\boldsymbol{\Sigma}^{-1} \frac{\partial \boldsymbol{\Sigma}}{\partial \alpha_i} \right) = -\text{Tr} \left(\boldsymbol{\Sigma} \frac{\partial \boldsymbol{\Sigma}^{-1}}{\partial \alpha_i} \right) = -\Sigma_{ii}, \quad \mathbf{m}^T \boldsymbol{\Sigma}^{-1} \frac{\partial \mathbf{m}}{\partial \alpha_i} =$$

$$-\mathbf{m}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma} \beta \frac{\partial \boldsymbol{\Sigma}^{-1}}{\partial \alpha_i} \boldsymbol{\Phi}^T \mathbf{t} = -m_i^2 \text{ and setting equal to 0:}$$

$$\frac{1}{2\alpha_i} + \frac{1}{2} m_i^2 - m_i^2 - \frac{1}{2} \Sigma_{ii} = 0 \rightarrow \alpha_i = \frac{1 - \alpha_i \Sigma_{ii}}{m_i^2} \equiv \frac{\gamma_i}{m_i^2}$$

RVM: Evidence Approximation

$$p(\mathcal{D}|\boldsymbol{\alpha}, \beta) = (2\pi)^{-\frac{N}{2}} \beta^{\frac{N}{2}} |\text{diag}(\boldsymbol{\alpha})|^{\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\beta \mathbf{t}^T \mathbf{t} - \mathbf{m}^T \boldsymbol{\Sigma}^{-1} \mathbf{m}) \right\} |\boldsymbol{\Sigma}|^{\frac{1}{2}}$$

□ Taking derivatives wrt β of $\ln p(\mathcal{D}|\boldsymbol{\alpha}, \beta)$ and using $\frac{\partial \ln |\boldsymbol{\Sigma}|}{\partial \beta} = \text{Tr} \left(\boldsymbol{\Sigma}^{-1} \frac{\partial \boldsymbol{\Sigma}}{\partial \beta} \right) = -\text{Tr} \left(\boldsymbol{\Sigma} \frac{\partial \boldsymbol{\Sigma}^{-1}}{\partial \beta} \right) = -\text{Tr}(\boldsymbol{\Sigma} \boldsymbol{\Phi}^T \boldsymbol{\Phi})$, and $\mathbf{m} = \beta \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \mathbf{t}$, setting equal to 0:

$$\frac{N}{2\beta} - \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{1}{2} \mathbf{m}^T \boldsymbol{\Phi}^T \boldsymbol{\Phi} \mathbf{m} + \mathbf{m}^T \boldsymbol{\Phi}^T \mathbf{t} - \frac{1}{2} \text{tr}(\boldsymbol{\Sigma} \boldsymbol{\Phi}^T \boldsymbol{\Phi}) = 0 \rightarrow (\beta^{(new)})^{-1} = \frac{\|\mathbf{t} - \boldsymbol{\Phi} \mathbf{m}\|^2}{N - \sum \gamma_i}$$

□ Here, we used

$$\begin{aligned} \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \boldsymbol{\Phi} &= \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \boldsymbol{\Phi} + \beta^{-1} \boldsymbol{\Sigma} \text{diag}(\boldsymbol{\alpha}) - \beta^{-1} \boldsymbol{\Sigma} \text{diag}(\boldsymbol{\alpha}) = \\ &= \boldsymbol{\Sigma} (\text{diag}(\boldsymbol{\alpha}) + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}) \beta^{-1} - \beta^{-1} \boldsymbol{\Sigma} \text{diag}(\boldsymbol{\alpha}) \\ &= (\text{diag}(\boldsymbol{\alpha}) + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} (\text{diag}(\boldsymbol{\alpha}) + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}) \beta^{-1} - \beta^{-1} \boldsymbol{\Sigma} \text{diag}(\boldsymbol{\alpha}) \\ &= \beta^{-1} (\mathbf{I} - \text{diag}(\boldsymbol{\alpha}) \boldsymbol{\Sigma}) \end{aligned}$$

and thus $\text{tr}(\boldsymbol{\Sigma} \boldsymbol{\Phi}^T \boldsymbol{\Phi}) = \beta^{-1} \sum \gamma_i$

RVM: Evidence Approximation

$$\alpha_i = \frac{1 - \alpha_i \Sigma_{ii}}{m_i^2} \equiv \frac{\gamma_i}{m_i^2}, \quad (\beta^{(new)})^{-1} = \frac{\|t - \Phi \mathbf{m}\|^2}{N - \sum \gamma_i}$$

- Learning therefore proceeds by choosing initial values for α and β , evaluating the mean and covariance of the posterior using $\Sigma = (\text{diag}(\alpha) + \beta \Phi^T \Phi)^{-1}$, $\mathbf{m} = \beta \Sigma \Phi^T \mathbf{t}$, respectively, and then alternately re-estimating the hyperparameters, using the Eqs. above and iterating until a suitable convergence criterion is satisfied.
- The second approach is to use the EM algorithm. These two approaches to finding the values of the hyperparameters that maximize the evidence are formally equivalent. Numerically, however, it is found that the direct optimization approach gives somewhat faster convergence (Tipping, 2001).

RVM: Fast training algorithm

- ❑ We will develop a **fast algorithm** for **maximizing the evidence**.
- ❑ We follow closely the developments in [*Tipping, M. E. and A. C. Faul \(2003\). Fast marginal likelihood maximisation for sparse Bayesian models. In C. M. Bishop and B. J. Frey \(Eds.\), Proceedings of the Ninth International Workshop on Artificial Intelligence and Statistics, Key West, FL, Jan 3-6.*](#)
- ❑ There are some typos [in the paper](#).
- ❑ We will achieve this by asking the following question:
- ❑ A Matlab implementation for the one dimensional case can be found [here](#).

Keeping β fixed and given some values for α , what is the action on a single basis function that gives the maximum increase in evidence?

This will require a very thorough examination of the evidence function.

RVM: Splitting the evidence

□ Some notation:

- Subscript “ $-i$ ”: Wherever we see this applied to a vector or matrix, it gives the corresponding vector or matrix with the influence of the i –th basis function removed.
- $\alpha_{-i} = (\alpha_1, \dots, \alpha_{i-1}, \alpha_{i+1}, \dots, \alpha_m)$: All the hyper-parameters except the ones that pertain to basis function i .
- Φ_{-i} : The design matrix with the column corresponding to basis function i removed.
- $\mathbf{C}_{-i} = \beta^{-1} \mathbf{I}_N + \Phi_{-i} \text{diag}(\alpha_{-i})^{-1} \Phi_{-i}^T$: The matrix \mathbf{C} , without the influence of basis function i .
- $\mathcal{E}(\alpha_{-i}, \beta) = -\frac{1}{2} [n \log(2\pi) + \log|\mathbf{C}_{-i}| + \mathbf{y}^T \mathbf{C}_{-i}^{-1} \mathbf{y}]$: The evidence when the i –th basis function has been removed.

RVM: Splitting the evidence

Given any basis function $i \in \{1, \dots, M\}$, the evidence can be split as:

$$\mathcal{E}(\boldsymbol{\alpha}, \beta) = \mathcal{E}(\boldsymbol{\alpha}_{-i}, \beta) + \epsilon(\alpha_i),$$

where

$$\epsilon(\alpha_i) = \frac{1}{2} \left(\log \alpha_i - \log(\alpha_i + s_i) + \frac{h_i^2}{\alpha_i + s_i} \right),$$

with

$$s_i = \boldsymbol{\varphi}_i^T \mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i \quad \text{and} \quad h_i = \boldsymbol{\varphi}_i^T \mathbf{C}_{-i}^{-1} \mathbf{t}.$$

Proof:

Remark:
Everything that
relates to basis
function i is there!

We start by splitting \mathbf{C} :

$$\begin{aligned} \mathbf{C} &= \beta^{-1} \mathbf{I}_N + \boldsymbol{\Phi} \text{diag}(\boldsymbol{\alpha})^{-1} \boldsymbol{\Phi}^T \\ &= \beta^{-1} \mathbf{I}_N + (\boldsymbol{\Phi}_{-i} \quad \boldsymbol{\varphi}_i) \begin{pmatrix} \text{diag}(\boldsymbol{\alpha}_{-i})^{-1} & \mathbf{0}_{(M-1) \times 1} \\ \mathbf{0}_{1 \times (M-1)} & \alpha_i^{-1} \end{pmatrix} \begin{pmatrix} \boldsymbol{\Phi}_{-i}^T \\ \boldsymbol{\varphi}_i^T \end{pmatrix} \\ &= \beta^{-1} \mathbf{I}_N + \boldsymbol{\Phi}_{-i} \text{diag}(\boldsymbol{\alpha}_{-i})^{-1} \boldsymbol{\Phi}_{-i}^T + \alpha_i^{-1} \boldsymbol{\varphi}_i \boldsymbol{\varphi}_i^T \\ &= \mathbf{C}_{-i} + \alpha_i^{-1} \boldsymbol{\varphi}_i \boldsymbol{\varphi}_i^T. \end{aligned}$$

RVM: Splitting the evidence

Proof (continuation):

Remember that the evidence is:

$$\mathcal{E}(\alpha, \beta) = -\frac{1}{2} [N \log(2\pi) + \log|\mathbf{C}| + \mathbf{t}^T \mathbf{C}^{-1} \mathbf{t}].$$

We will only make use of this:

$$\mathbf{C} = \mathbf{C}_{-i} + \alpha_i^{-1} \boldsymbol{\varphi}_i \boldsymbol{\varphi}_i^T,$$

and some identities.

From the Woodbury identity, we get:

$$\begin{aligned} & (\mathbf{A} + \mathbf{UCV})^{-1} \\ &= \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{U} (\mathbf{C}^{-1} + \mathbf{VA}^{-1} \mathbf{U})^{-1} \mathbf{VA}^{-1}. \end{aligned}$$

$$\begin{aligned} \mathbf{C}^{-1} &= \mathbf{C}_{-i}^{-1} - \mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i (\alpha_i + \boldsymbol{\varphi}_i^T \mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i)^{-1} \boldsymbol{\varphi}_i^T \mathbf{C}_{-i}^{-1} \\ &= \mathbf{C}_{-i}^{-1} - \frac{\mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i \boldsymbol{\varphi}_i^T \mathbf{C}_{-i}^{-1}}{\alpha_i + s_i}. \end{aligned}$$

And we get:

$$\mathbf{t}^T \mathbf{C}^{-1} \mathbf{t} = \mathbf{t}^T \mathbf{C}_{-i}^{-1} \mathbf{t} - \frac{(\mathbf{t}^T \mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i)(\boldsymbol{\varphi}_i^T \mathbf{C}_{-i}^{-1} \mathbf{t})}{\alpha_i + s_i} = \mathbf{t}^T \mathbf{C}_{-i}^{-1} \mathbf{t} + \frac{h_i^2}{\alpha_i + s_i}.$$

RVM: Splitting the evidence

Proof (continuation):

Finally, we look at $|\mathbf{C}|$. Remember:

$$\mathbf{C} = \mathbf{C}_{-i} + \alpha_i^{-1} \boldsymbol{\varphi}_i \boldsymbol{\varphi}_i^T,$$

Using the matrix determinant lemma:

$$|\mathbf{A} + \mathbf{U}\mathbf{W}\mathbf{V}^T| = |\mathbf{W}^{-1} + \mathbf{V}^T \mathbf{A}^{-1} \mathbf{U}| |\mathbf{W}| |\mathbf{A}|$$

$$\begin{aligned} |\mathbf{C}| &= |\alpha_i + \boldsymbol{\varphi}_i^T \mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i| |\alpha_i^{-1}| |\mathbf{C}_{-i}| \\ &= \frac{|\mathbf{C}_{-i}| (\alpha_i + s_i)}{\alpha_i}. \end{aligned}$$

Or

$$\log |\mathbf{C}| = \log |\mathbf{C}_{-i}| - \log \alpha_i + \log(\alpha_i + s_i).$$

RVM: Splitting the evidence

Proof (continuation):

Now we combine everything. Remember that:

$$\mathbf{t}^T \mathbf{C}^{-1} \mathbf{t} = \mathbf{t}^T \mathbf{C}_{-i}^{-1} \mathbf{t} + \frac{h_i^2}{\alpha_i + s_i},$$

and

$$\log |\mathbf{C}| = \log |\mathbf{C}_{-i}| - \log \alpha_i + \log(\alpha_i + s_i).$$

Therefore, we have for the evidence:

$$\begin{aligned} \mathcal{E}(\boldsymbol{\alpha}, \beta) &= -\frac{1}{2} [N \log(2\pi) + \log |\mathbf{C}| + \mathbf{t}^T \mathbf{C}^{-1} \mathbf{t}] \\ &= -\frac{1}{2} \left(N \log(2\pi) + \log |\mathbf{C}_{-i}| + \mathbf{t}^T \mathbf{C}_{-i}^{-1} \mathbf{t} + \frac{h_i^2}{\alpha_i + s_i} - \log \alpha_i + \log(\alpha_i + s_i) \right) \\ &= \mathcal{E}(\boldsymbol{\alpha}_{-i}, \beta) + \epsilon(\alpha_i). \end{aligned}$$

End of proof.

RVM: Analysis of the evidence

We have shown that for any basis function $i \in \{1, \dots, M\}$, the evidence can be split as:

$$\mathcal{E}(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \mathcal{E}(\boldsymbol{\alpha}_{-i}, \boldsymbol{\beta}) + \epsilon(\alpha_i),$$

where

$$\epsilon(\alpha_i) = \frac{1}{2} \left(\log \alpha_i - \log(\alpha_i + s_i) + \frac{h_i^2}{\alpha_i + s_i} \right).$$

We will now analyze the $\epsilon(\alpha_i)$ term in more detail.

We compute its first derivative:

$$\frac{\partial \epsilon(\alpha_i)}{\partial \alpha_i} = \frac{s_i^2 - \alpha_i \theta_i}{2\alpha_i(\alpha_i + s_i)^2},$$

where

$$\theta_i = h_i^2 - s_i.$$

RVM: Analysis of the evidence

Setting

$$\frac{\partial \epsilon(\alpha_i)}{\partial \alpha_i} = \frac{s_i^2 - \alpha_i \theta_i}{2\alpha_i(\alpha_i + s_i)^2}$$

equal to zero, we have two possibilities:

□ If $\theta_i > 0$, then

$$\alpha_i^* = \frac{s_i^2}{\theta_i},$$

is the unique maximum of $\epsilon(\alpha_i)$ (calculate the 2nd derivative to be convinced that it is indeed a maximum).

□ If $\theta_i \leq 0$, then

$$\alpha_i^* = +\infty,$$

Is the unique maximum of $\epsilon(\alpha_i)$ (convince yourselves that its derivative remains positive for all α_i). This means that removing the corresponding basis function from the model, increases the evidence.

RVM: Possible actions

The previous observations translate into three possible actions for each basis function i :

□ **ADD**, if $\theta_i > 0$ and $\alpha_i = +\infty$:

$$\alpha_i^{\text{new}} = \frac{s_i^2}{\theta_i},$$
$$\Delta\mathcal{E} = \epsilon(\alpha_i^{\text{new}}).$$

□ **RE-ESTIMATE**, if $\theta_i > 0$ and $\alpha_i < +\infty$:

$$\alpha_i^{\text{new}} = \frac{s_i^2}{\theta_i},$$
$$\Delta\mathcal{E} = \epsilon(\alpha_i^{\text{new}}) - \epsilon(\alpha_i).$$

□ **DELETE**, if $\theta_i \leq 0$ and $\alpha_i < +\infty$:

$$\alpha_i^{\text{new}} = +\infty,$$
$$\Delta\mathcal{E} = -\epsilon(\alpha_i).$$

RVM: Updating the noise

Now, we assume that all the α 's are kept fixed and we attempt to maximize the evidence with respect to the inverse noise β .

Taking the derivative of the evidence with respect to β , you get:

$$\beta^{\text{new}} = \frac{N - M - \sum_i \alpha_i \Sigma_{ii}}{\| \mathbf{t} - \mathbf{\Phi} \boldsymbol{\mu} \|^2}.$$

This is a lengthy calculation. You are going to need the following formulas for a general matrix \mathbf{A} :

$$\frac{d\mathbf{A}^{-1}}{dt} = -\mathbf{A}^{-1} \frac{d\mathbf{A}}{dt} \mathbf{A}^{-1},$$

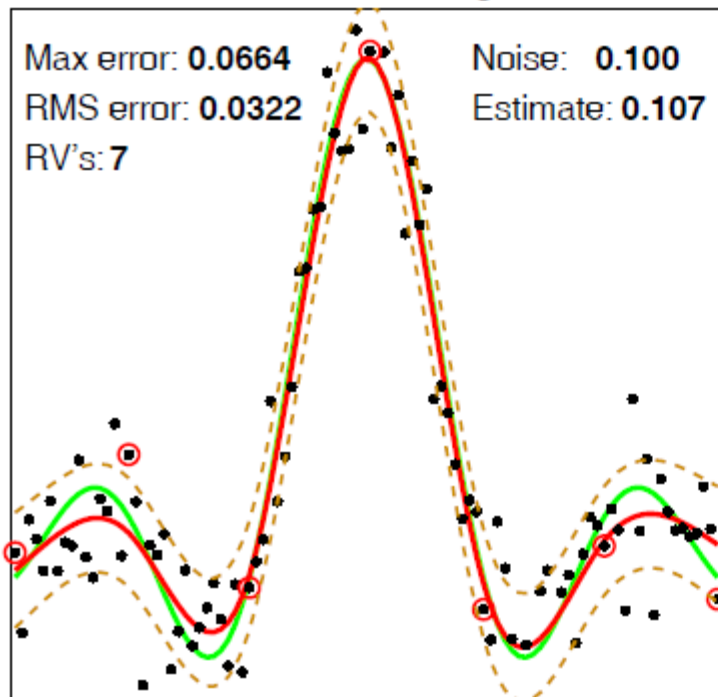
$$\frac{d|\mathbf{A}|}{dt} = |\mathbf{A}| \text{tr} \left[\mathbf{A}^{-1} \frac{d\mathbf{A}}{dt} \right].$$

RVM: Complete Algorithm

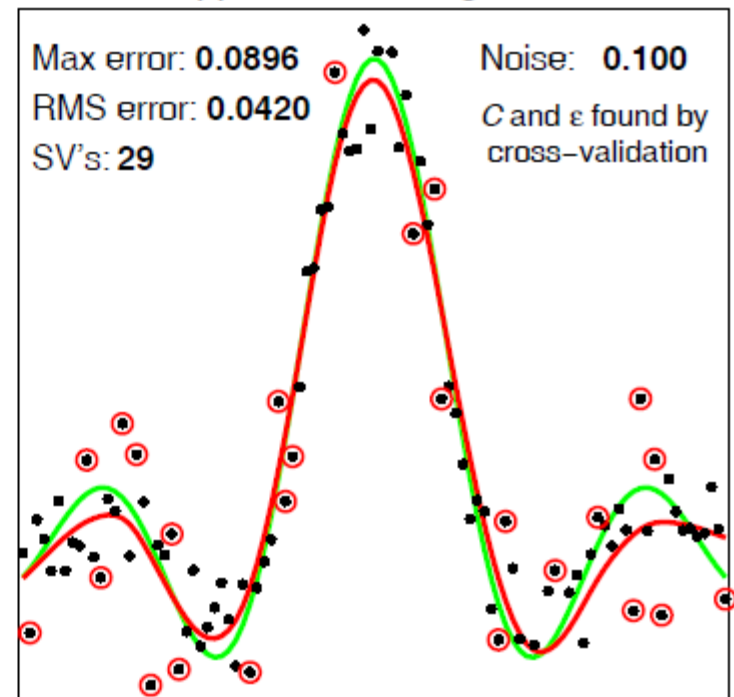
1. Initialize β and α (Some of them might be $+\infty$. The corresponding basis functions are not in the model).
2. Compute: Σ , \mathbf{m} only for the basis functions in the model but **all** the s_i and h_i 's.
3. For $i = 1, \dots, M$, compute $\theta_i = h_i^2 - s_i$ and find the basis function i^* and the corresponding action \mathcal{A}_{i^*} that gives the maximum increase in evidence.
4. Perform the best action.
5. Re-estimate the noise β .
6. Re-compute (or update) Σ , \mathbf{m} , s_i and h_i .
7. If converged, stop. Otherwise go to 3.

A simple example

Relevance Vector Regression



Support Vector Regression



Fitting noisy samples from $f(x) = \frac{\sin(x)}{x}$. The basis functions are Gaussian kernels centered on the data points:

$$\phi_i(x) = k(x, x^{(i)}) = \exp\left(-\lambda(x - x^{(i)})^2\right).$$

Taken [from this](#) very nice [tutorial](#). The comparison is with [SVM](#) (a non-Bayesian competitor).

RVM: Fast updates

An efficient implementation of the RVM algorithm, requires a fast way to update the various statistics.

First, the mean and the co-variance matrix should only be build out of the **relevant** basis functions only:

$$\begin{aligned}\Sigma_r &= (\text{diag}(\alpha_r) + \beta \Phi_r^T \Phi_r)^{-1}, \\ \mathbf{m}_r &= \beta \Sigma_r \Phi_r^T \mathbf{y},\end{aligned}$$

where the subscript “r” denotes that the corresponding quantity corresponds only to the basis functions in the model (the relevant basis function, i.e. the ones with $\alpha_i < +\infty$. The other basis functions are simply irrelevant...

This requires some book-keeping. That is, at any iteration you must keep track of the set of relevant basis functions:

$$\mathcal{R} = \{j \in \{1, \dots, m\}: \alpha_j < +\infty\}.$$

RVM: Fast updates

Secondly, we come to the statistics:

$$s_i = \boldsymbol{\varphi}_i^T \mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i \quad \text{and} \quad h_i = \boldsymbol{\varphi}_i^T \mathbf{C}_{-i}^{-1} \mathbf{t}.$$

These are required for **all** basis functions, so that the change in evidence can be computed. It is more convenient to keep track of the following statistics instead:

$$S_i = \boldsymbol{\varphi}_i^T \mathbf{C}^{-1} \boldsymbol{\varphi}_i \quad \text{and} \quad H_i = \boldsymbol{\varphi}_i^T \mathbf{C}^{-1} \mathbf{t}.$$

It is trivial to show that they connect with the original by:

$$s_i = \frac{\alpha_i S_i}{\alpha_i - S_i} \quad \text{and} \quad h_i = \frac{\alpha_i H_i}{\alpha_i - S_i}.$$

Using the Woodbury identity once more, you can find an easy way to compute the capital statistics:

$$\begin{aligned} S_i &= \beta \|\boldsymbol{\varphi}_i\|^2 - \beta^2 \boldsymbol{\varphi}_i^T \boldsymbol{\Phi}_r \boldsymbol{\Sigma}_r \boldsymbol{\Phi}_r^T \boldsymbol{\varphi}_i, \\ H_i &= \beta \boldsymbol{\varphi}_i^T \mathbf{t} - \beta^2 \boldsymbol{\varphi}_i^T \boldsymbol{\Phi}_r \boldsymbol{\Sigma}_r \boldsymbol{\Phi}_r^T \mathbf{t}. \end{aligned}$$

RVM: Fast updates

Finally, each action can be carried out without computing the statistics from scratch! You can derive them by repeatedly applying the Woodbury identity. It will take a couple of pages...

Notation:

- M_r : current number of relevant basis functions
- $i \in \{1, \dots, M\}$: index of basis function that is currently being updated.
- $j \in \{1, \dots, M_r\}$: the index within the **relevant** basis set that corresponds to i .
- $t \in \{1, \dots, M\}$: index that runs over all basis functions (in or out of the model).

RVM: Fast updates

ADD basis function i :

$$\begin{aligned}\Delta\mathcal{E}^{\text{new}} &= \frac{1}{2} \left(\frac{Q_i^2 - S_i}{S_i} + \log \frac{S_i}{Q_i^2} \right), \\ \Sigma_r^{\text{new}} &= \begin{pmatrix} \Sigma_r + \beta^2 \Sigma_{ii} \Sigma_r \Phi_r^T \varphi_i \varphi_i^T \Phi_r \Sigma_r & -\beta \Sigma_{ii} \Sigma_r \Phi_r^T \varphi_i \\ -\beta \Sigma_{ii} \varphi_i^T \Phi_r \Sigma_r & \Sigma_{ii} \end{pmatrix}, \\ \mathbf{m}_r^{\text{new}} &= \begin{pmatrix} \mathbf{m}_r - \mu_i \beta \Sigma_{ii} \Sigma_r \Phi_r^T \varphi_i \\ m_i \end{pmatrix}, \\ S_t^{\text{new}} &= S_t - \Sigma_{ii} (\beta \varphi_t^T \mathbf{e}_i)^2, \\ H_t^{\text{new}} &= H_t - \mu_i \beta \varphi_t^T \mathbf{e}_i,\end{aligned}$$

where

$$\begin{aligned}\Sigma_{ii} &= (\alpha_i^{\text{new}} - S_i)^{-1}, \\ \mu_i &= \Sigma_{ii} H_i,\end{aligned}$$

and

$$\mathbf{e}_i = \varphi_i - \beta \Phi_r \Sigma_r \Phi_r^T \varphi_i.$$

Finally, m_r is increased by one.

CAUTION: There is a typo in Tipping and Faul.

RVM: Fast updates

DELETE basis function j :

$$\begin{aligned}\Delta\mathcal{E}^{\text{new}} &= \frac{1}{2} \left(\frac{Q_i^2 - S_i}{S_i} - \log \left(1 - \frac{S_i}{\alpha_i} \right) \right), \\ \Sigma_r^{\text{new}} &= \Sigma_r - \Sigma_{jj}^{-1} \Sigma_j \Sigma_j^T, \\ \mathbf{m}_r^{\text{new}} &= \mathbf{m}_r - \mu_j \Sigma_{jj}^{-1} \Sigma_j, \\ S_t^{\text{new}} &= S_t + \Sigma_{jj}^{-1} (\beta \Sigma_j^T \Phi_r^T \boldsymbol{\varphi}_t)^2, \\ H_t^{\text{new}} &= H_t + m_j \Sigma_{jj}^{-1} \beta \Sigma_j^T \Phi_r^T \boldsymbol{\varphi}_t,\end{aligned}$$

where Σ_{jj} is the (j, j) element and Σ_j the j -th column of Σ_r , resp., and m_j is the j -th element of \mathbf{m}_r .

After completing this step the corresponding rows and columns of the mean and the covariance must be removed.

The indices of the relevant vectors must be updated accordingly.

RVM: Fast updates

RE-ESTIMATE basis function j :

$$\Delta \mathcal{E}^{\text{new}} = \frac{1}{2} \left(\frac{Q_i^2}{S_i + \left((\alpha_i^{\text{new}})^{-1} - \alpha_i^{-1} \right)^{-1}} - \log \left(1 - S_i \left((\alpha_i^{\text{new}})^{-1} - \alpha_i^{-1} \right) \right) \right),$$

$$\mathbf{\Sigma}_r^{\text{new}} = \mathbf{\Sigma}_r - \kappa_j \mathbf{\Sigma}_j \mathbf{\Sigma}_j^T,$$

$$\mathbf{m}_r^{\text{new}} = \mathbf{m}_r - \kappa_j m_j \mathbf{\Sigma}_j,$$

$$S_t^{\text{new}} = S_t + \kappa_j \left(\beta \mathbf{\Sigma}_j^T \mathbf{\Phi}_r^T \boldsymbol{\varphi}_t \right)^2,$$

$$H_t^{\text{new}} = H_t + \kappa_j m_j \beta \mathbf{\Sigma}_j^T \mathbf{\Phi}_r^T \boldsymbol{\varphi}_t,$$

$\mathbf{\Sigma}_j$ the j -th column of $\mathbf{\Sigma}_r$, resp., and μ_j is the j -th element of \mathbf{m}_r and

$$\kappa_j = \left(\Sigma_{jj} + \left((\alpha_i^{\text{new}})^{-1} - \alpha_i^{-1} \right)^{-1} \right)^{-1},$$

where Σ_{jj} is the (j, j) element.

RVM: Extension to multiple outputs

- We assume that we observe the data:

$$\mathcal{D} = \{\mathbf{x}^{(i)}, \mathbf{t}^{(i)} = \mathbf{f}(\mathbf{x}^{(i)})\}_{i=1}^n,$$

and that we wish to learn the function $\mathbf{f}(\cdot) \in \mathbf{R}^q$.

- Pick a set of M basis functions $\boldsymbol{\phi}: \mathbf{R}^k \rightarrow \mathbf{R}^M$:

$$\boldsymbol{\phi}(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_M(\mathbf{x})).$$

- We will approximate $\mathbf{f}(\cdot)$ by:

$$\hat{\mathbf{f}}(\mathbf{x}; \mathbf{W}) = \mathbf{W}^T \boldsymbol{\phi}(\mathbf{x}) + \boldsymbol{\epsilon},$$

where $\boldsymbol{\epsilon}$ is a random vector representing noise and $\mathbf{W} \in \mathbf{R}^{M \times q}$.

- For simplicity, we choose to work with Gaussian noise with zero mean and inverse variance β , so the **likelihood** is:

$$p(\mathbf{t}|\mathbf{x}, \mathbf{W}, \beta) = \mathcal{N}_q(\mathbf{t}|\mathbf{W}^T \boldsymbol{\phi}(\mathbf{x}), \beta^{-1}).$$

RVM: Extension to multiple outputs

- The prior on the weights:

$$p(\mathbf{W}|\boldsymbol{\alpha}) = \prod_{j=1}^M p(\mathbf{w}_j|\alpha_j),$$

with $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_M)$ and:

$$p(\mathbf{w}_j|\alpha_j) = \mathcal{N}_q(\mathbf{w}_j|0, \alpha_j^{-1}).$$

- The posterior of the weights is:

$$p(\mathbf{W}|\mathcal{D}, \beta, \boldsymbol{\alpha}) = \mathcal{N}(\mathbf{W}|\mathbf{M}, \boldsymbol{\Sigma}),$$

where:

$$\begin{aligned}\boldsymbol{\Sigma} &= (\text{diag}(\boldsymbol{\alpha}) + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1}, \\ \mathbf{M} &= \beta \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \mathbf{Y},\end{aligned}$$

with $\mathbf{Y} = (\mathbf{t}^{(1)}, \dots, \mathbf{t}^{(N)}) \in \mathbf{R}^{N \times q}$ is the vector of observations, $\boldsymbol{\Phi} \in \mathbf{R}^{N \times M}$ is the design matrix:

RVM: Predictive distribution

□ The **predictive** distribution is easily found as follows:

$$\begin{aligned} p(\mathbf{t}|\mathbf{x}, \alpha, \beta) &= \int p(\mathbf{t}|\mathbf{x}, \mathbf{W}, \beta) p(\mathbf{W}|\alpha) d\mathbf{W} \\ &= \mathcal{N}(y|\mathbf{M}^T \boldsymbol{\phi}(\mathbf{x}), \boldsymbol{\varphi}(\mathbf{x})^T \boldsymbol{\Sigma} \boldsymbol{\varphi}(\mathbf{x}) + \beta^{-1}). \end{aligned}$$

So, we have the *predictive mean*:

$$\mu(\mathbf{x}) = \mathbf{M}^T \boldsymbol{\phi}(\mathbf{x}),$$

and the predictive variance:

$$\sigma^2(\mathbf{x}) = \boldsymbol{\varphi}(\mathbf{x})^T \boldsymbol{\Sigma} \boldsymbol{\varphi}(\mathbf{x}) + \beta^{-1}.$$

Proof:

Complete the square....

End of proof.

RVM: Extension to multiple outputs

- **REMARK:** You might have to work with a scaled version of the data so that the prior assumptions are closer to reality. You may attempt to fit the function:

$$\mathbf{g}(\mathbf{x}) = \mathbf{L}_{\text{obs}}^{-1}(\mathbf{f}(\mathbf{x}) - \mathbf{m}_{\text{obs}}),$$

where \mathbf{m}_{obs} is the *empirical mean of the data*:

$$\mathbf{m}_{\text{obs}} = \frac{1}{n} \sum_{i=1}^N \mathbf{t}^{(i)},$$

and $\mathbf{L}_{\text{obs}} \in \mathbf{R}^{q \times q}$ is the Cholesky decomposition of the *empirical covariance matrix of the data*:

$$\mathbf{C}_{\text{obs}} = \frac{1}{N-1} \sum_{i=1}^N (\mathbf{t}^{(i)} - \mathbf{m}_{\text{obs}})(\mathbf{t}^{(i)} - \mathbf{m}_{\text{obs}})^T.$$

If the empirical covariance is rank-deficient, you may use either an incomplete Cholesky or the eigen-decomposition. This is essentially a sort of output dimensionality reduction.

RVM: Extension to multiple outputs

- We proceed again, by maximizing the total evidence.
- Working out the details, it is easy to show that the total evidence is the sum of the evidence of each output:

$$\mathcal{E}(\boldsymbol{\alpha}, \beta | \mathcal{D}) = \sum_{j=1}^q \mathcal{E}(\boldsymbol{\alpha}, \beta | \mathcal{D}_j),$$

where $\mathcal{D}_j = \left\{ \left(\mathbf{x}^{(i)}, t_j^{(i)} \right) \right\}_{i=1}^N$.

- The algorithmic details remain essentially the same. The only difference is that:

$$H_i = \frac{1}{q} \sum_{j=1}^q H_{ji},$$

$$H_{ji} = \beta \boldsymbol{\varphi}_i^T \mathbf{t}_j - \beta^2 \boldsymbol{\varphi}_i^T \boldsymbol{\Phi}_r \boldsymbol{\Sigma}_r \boldsymbol{\Phi}_r^T \mathbf{t}_j.$$

RVM: Numerical stability concerns

- ❑ There are several numerical instabilities that should concern you.
- ❑ This is very common.
- ❑ For a numerically stable way to implement RVM using the Generalized Singular Value Decomposition look at Appendix D of [I. Bilonis and N. Zabarar. *Multidimensional adaptive relevance vector machines for uncertainty quantification*. SIAM Journal for Scientific Computing, Vol. 34, No. 6, pp. B881–B908 2012.](#)

RVM For Classification

- We can extend the relevance vector machine framework to classification problems by applying the ARD prior over weights to a probabilistic linear classification model.
- Consider two-class problems with a binary target variable $t \in \{0, 1\}$.
- The model now takes the form of a linear combination of basis functions transformed by a logistic sigmoid function

$$y(\mathbf{x}, \mathbf{w}) = \sigma(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}))$$

- If we introduce a Gaussian prior over \mathbf{w} , then we obtain the model that has been considered in the logistic regression notes.

RVM For Classification

- In the RVM, this model uses the ARD prior

$$p(w_j | \alpha_j) = \mathcal{N}(w_j | 0, \alpha_j^{-1}).$$

- We can no longer integrate analytically over \mathbf{w} . Here we follow Tipping (2001) and use the Laplace approximation.
- We begin by initializing $\boldsymbol{\alpha}$. For this given value of $\boldsymbol{\alpha}$, we then build a Gaussian approximation to the posterior distribution and thereby obtain an approximation to the marginal likelihood ($A = \text{diag}(\alpha_i)$)

$$\begin{aligned} \ln p(\mathbf{w} | \mathbf{t}, \boldsymbol{\alpha}) &= \ln \{p(\mathbf{t} | \mathbf{w}) p(\mathbf{w} | \boldsymbol{\alpha})\} - \ln p(\mathbf{t} | \boldsymbol{\alpha}) \\ &= \sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} - \frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} + \text{const.} \end{aligned}$$

- Maximization of this approximate marginal likelihood then leads to a re-estimated value for $\boldsymbol{\alpha}$, and the process is repeated until convergence.

RVM For Classification

- This can be done using iterative reweighted least squares (IRLS).
- We need the gradient and Hessian of the log posterior distribution

$$\begin{aligned}\nabla \ln p(\mathbf{w}|\mathbf{t}, \alpha) &= \mathbf{\Phi}^T(\mathbf{t} - \mathbf{y}) - \mathbf{A}\mathbf{w} \\ \nabla \nabla \ln p(\mathbf{w}|\mathbf{t}, \alpha) &= -(\mathbf{\Phi}^T \mathbf{B} \mathbf{\Phi} + \mathbf{A})\end{aligned}$$

- \mathbf{B} is an $N \times N$ diagonal matrix with elements $b_n = y_n(1 - y_n)$, the vector $\mathbf{y} = (y_1, \dots, y_N)^T$, and $\mathbf{\Phi}$ is the design matrix with elements $\Phi_{ni} = \phi_i(\mathbf{x}_n)$.
- At convergence of the IRLS algorithm, the negative Hessian represents the inverse covariance matrix for the Gaussian approximation to the posterior distribution.
- The mean and covariance of the Laplace approximation in the form

$$\begin{aligned}\mathbf{w} &= \mathbf{A}^{-1} \mathbf{\Phi}^T(\mathbf{t} - \mathbf{y}) \\ \mathbf{\Sigma} &= (\mathbf{\Phi}^T \mathbf{B} \mathbf{\Phi} + \mathbf{A})^{-1}\end{aligned}$$

RVM For Classification

- The marginal likelihood can be approximated as:

$$p(\mathbf{t}|\boldsymbol{\alpha}) = \int p(\mathbf{t}|\mathbf{w})p(\mathbf{w}|\boldsymbol{\alpha})d\mathbf{w} \approx p(\mathbf{t}|\mathbf{w}^*)p(\mathbf{w}^*|\boldsymbol{\alpha})(2\pi)^{M/2}|\boldsymbol{\Sigma}|^{1/2}$$

- If we substitute for $p(\mathbf{t}|\mathbf{w}^*)$ and $p(\mathbf{w}^*|\boldsymbol{\alpha})$ and then set the derivative wrt equal to zero:

$$-\frac{1}{2}(w_i^*)^2 + \frac{1}{2\alpha_i} - \frac{1}{2}\Sigma_{ii} = 0$$

- Defining $\gamma_i = 1 - \alpha_i\Sigma_{ii}$, gives:

$$\alpha_i^{new} = \frac{\gamma_i}{(w_i^*)^2}$$

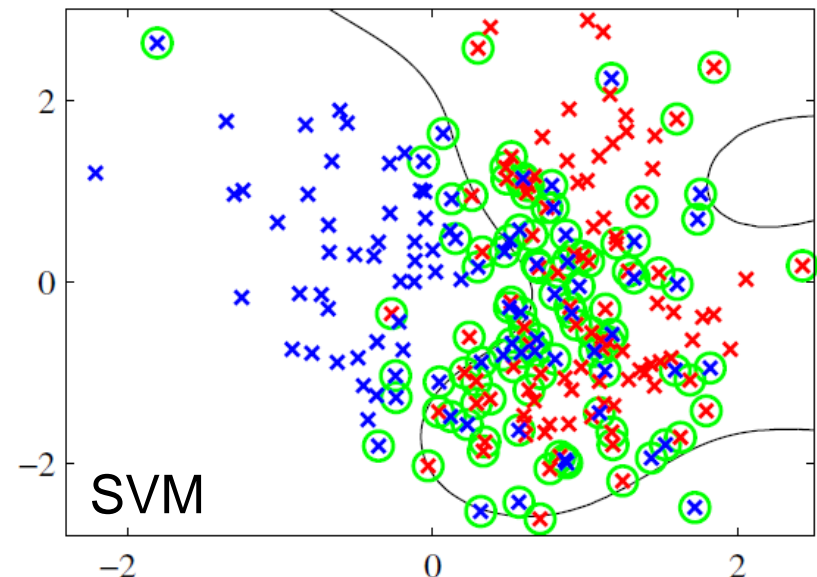
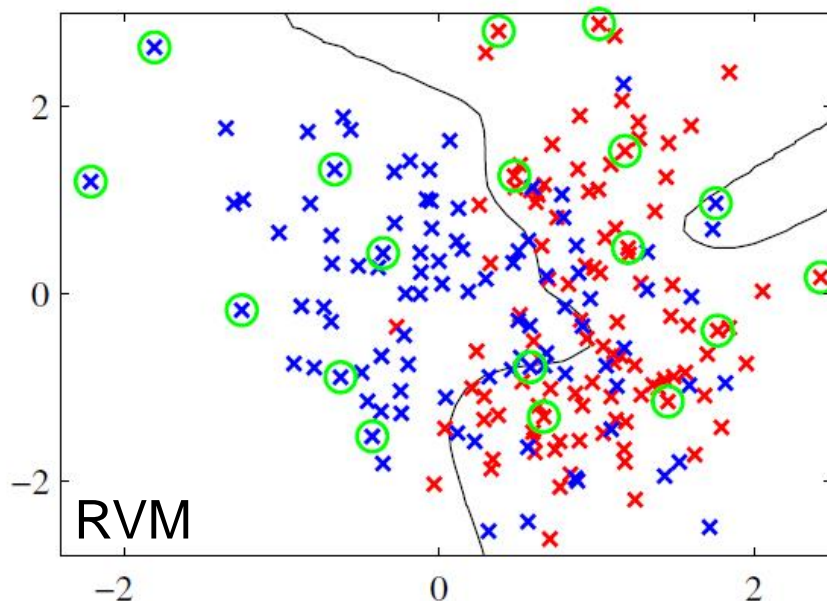
- With $\hat{\mathbf{t}} = \boldsymbol{\Phi}\mathbf{w}^* + \mathbf{B}^{-1}(\mathbf{t} - \mathbf{y})$, we can write

$$\ln p(\mathbf{t}|\boldsymbol{\alpha}) = -\frac{1}{2}\{N\ln(2\pi) + \ln|\mathbf{C}| + (\hat{\mathbf{t}})^T \mathbf{C}^{-1}\hat{\mathbf{t}}\}, \mathbf{C} = \mathbf{B} + \boldsymbol{\Phi}\mathbf{A}\boldsymbol{\Phi}^T$$

- This is the same form as for regression. Can apply the same sparsity analysis and obtain the same fast learning algorithm in which we optimize a single α_i at each step.

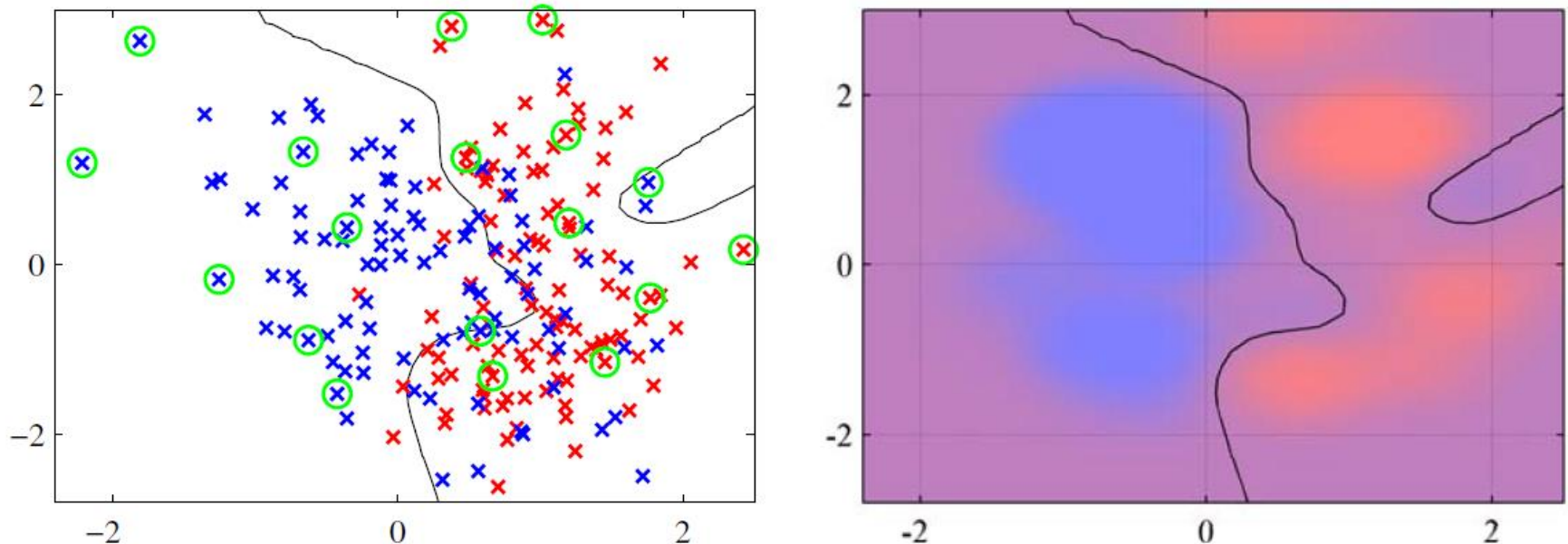
RVM For Classification

- Example of the RVM applied to a synthetic data set, in which the left-hand plot shows the decision boundary and the data points, with the relevance vectors indicated by circles.
- Comparing with the SVM results it is seen that the RVM gives a much sparser model.



RVM For Classification

- The right-hand plot shows the posterior probability given by the RVM output in which the proportion of red (blue) ink indicates the probability of that point belonging to the red (blue) class.



RVM For Classification

- So far, we have considered the RVM for binary classification problems. For $K > 2$ classes,

$$y_k(\mathbf{x}) = \frac{\exp(\alpha_k)}{\sum_j \exp(\alpha_k)}$$

- The log likelihood function is then given by

$$\ln p(\mathbf{T} | \mathbf{w}_1, \dots, \mathbf{w}_K) = \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}}$$

where the target values t_{nk} have a 1-of- K coding for each data point n , and \mathbf{T} is a matrix with elements t_{nk} . Again, the Laplace approximation can be used to optimize the hyperparameters in which the model and its Hessian are found using IRLS.

RVM For Classification

- ❑ The principal disadvantage is that the Hessian matrix has size $MK \times MK$, where M is the number of active basis functions, which gives an additional factor of K^3 in the computational cost of training compared with the two-class RVM.
- ❑ The principal disadvantage of the relevance vector machine is the relatively long training times compared with the SVM. This is offset, however, by the avoidance of cross-validation runs to set the model complexity parameters.
- ❑ Furthermore, because it yields sparser models, the computation time on test points, which is usually the more important consideration in practice, is typically much less.

References

- T. Fletcher, [Relevance Vector Machines Explained](#), October 2010.
- Tipping, M. (2001). [Sparse Bayesian learning and the relevance vector machine](#). *J. of Machine Learning Research* 1, 211–244.
- MacKay, D. (1995b). [Probable networks and plausible predictions — a review of practical Bayesian methods for supervised neural networks](#). *Network*.
- Neal, R. (1996). [Bayesian learning for neural networks](#). Springer.
- Bishop, C. and M. Tipping (2000). [Variational relevance vector machines](#). In *UAI*.
- Buntine, W. and A. Weigend (1991). [Bayesian backpropagation](#). *Complex Systems* 5, 603–643
- Tipping, M. (2001). [Sparse Bayesian learning and the relevance vector machine](#). *J. of Machine Learning Research* 1, 211–244.
- Wipf, D. and S. Nagarajan (2010, April). [Iterative Reweighted \$\ell_1\$ and \$\ell_2\$ Methods for Finding Sparse Solutions](#). *J. of Selected Topics in Signal Processing (Special Issue on Compressive Sensing)* 4(2).
- Wipf, D. and S. Nagarajan (2010, April). [Iterative Reweighted \$\ell_1\$ and \$\ell_2\$ Methods for Finding Sparse Solutions](#). *J. of Selected Topics in Signal Processing (Special Issue on Compressive Sensing)* 4(2).
- Wipf, D. and S. Nagarajan (2007). [A new view of automatic relevancy determination](#). In *NIPS*.
- MacKay, D. (1999). [Comparison of approximate methods for handling hyperparameters](#). *Neural Computation* 11(5), 1035–1068.
- Wipf, D. and S. Nagarajan (2007). [A new view of automatic relevancy determination](#). In *NIPS*.
- Wipf, D. and S. Nagarajan (2010, April). [Iterative Reweighted \$\ell_1\$ and \$\ell_2\$ Methods for Finding Sparse Solutions](#). *J. of Selected Topics in Signal Processing (Special Issue on Compressive Sensing)* 4(2).

References

- Wipf, D. and S. Nagarajan (2007). A new view of automatic relevancy determination. In *NIPS*.
- Rasmussen, C. E. and J. Quiñonero-Candela (2005). [Healing the relevance vector machine by augmentation](#). In L. D. Raedt and S. Wrobel (Eds.), *Proceedings of the 22nd International Conference on Machine Learning*, pp. 689–696.
- T. Fletcher, [Relevance Vector Machines Explained](#), October 2010.
- Faul, A. C. and M. E. Tipping (2002). Analysis of sparse Bayesian learning. In T. G. Dietterich, S. Becker, and Z. Ghahramani (Eds.), *Advances in Neural Information Processing Systems*, Volume 14, pp. 383–389. MIT Press.
- Tipping, M. E. and A. Faul (2003). Fast marginal likelihood maximization for sparse Bayesian models. In C. M. Bishop and B. Frey (Eds.), *Proceedings Ninth International Workshop on Artificial Intelligence and Statistics*, Key West, Florida.
- Williams, O., A. Blake, and R. Cipolla (2005). Sparse Bayesian learning for efficient visual tracking. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **27**(8), 1292–1304.